

Explanation of Mathematical Methods Behind the First Half of the Project

Emma Bowen

Introduction

What mathematical methods will I explain in this document?

- Gaussian processes
- Bayesian History Matching
- Kriging
- Expected Improvement
- Augmented Expected Improvement
- Knowledge Gradient

Context

We are building on the paper “Using history matching to speed up management strategy evaluation grid searches”. This paper is looking to find the Harvest Control Rule parametrised by F_{target} and $B_{trigger}$ that maximises the catch whilst keeping the risk below 0.05 for a single-stock fishery. The paper does not consider fleet dynamics. We have an objective function in the paper which combines the risk and the catch and so this is the function we want to maximise to get our maximal catch with the constraint of keeping the risk below 0.05.

To do this, the paper takes a Bayesian History Matching (BHM) approach. Firstly, we sample our first eight points which are spaced evenly throughout the sample space to get some initial data. Then, for each round we do the following:

- We set up or update the Gaussian Process (GP) to model the risk and the GP to model the catch

- We can then use the risk as a threshold so that we only consider the values of F_{target} and $B_{trigger}$ that have risk below 0.05
- We use the GP which is modelling the catch to get the value for the catch at every point in the sample space which will have some uncertainty
- We use BHM to remove any points that are implausible (that have a low probability of being higher than the current best catch)

We repeat this process until there is only one plausible point left and then we will accept this as being the F_{target} and $B_{trigger}$ that maximise the catch whilst keeping the risk below 0.05.

Now, we will look at the mathematics behind the methods above and also at the acquisition functions of Expected Improvement, Augmented Expected Improvement and Knowledge Gradient which I added to the original code.

Gaussian Processes

A Gaussian process F has a mean function μ_0 and a covariance matrix $\Sigma_{1:k}$ and is a probability distribution over our objective function f with the property that, for any given collection of points x_1, \dots, x_m , the marginal probability distribution on $(F(x_1), \dots, F(x_m))$ is given by(Frazier 2018):

$$(F(x_1), \dots, F(x_m)) \sim N((\mu_0(x_1), \dots, \mu_0(x_m)), \text{Cov}_{1:k}) \quad (1)$$

where $\Sigma_{1:k}$ can be found by evaluating the covariance function $\text{cov}_0(x_i, x_j)$ for every pair x_i, x_j where $j, i \in \{1, \dots, k\}$.

We choose a covariance function such that inputs that have nearby points that have been evaluated have a more certain output than points that are further away from the points that have been evaluated(Spence 2025). This is equivalent to saying that if for some x, x', x'' in the design space we have $\|x - x'\| < \|x - x''\|$ for some norm $\|\cdot\|$, then $\text{cov}_0(x, x') > \text{cov}_0(x, x'')$ (Frazier 2018).

We can set up GPs in R using the DiceKriging package as below:

```
gp_cat <- km(~.^2, design=runs[,c("Ftarget", "Btrigger")], estim.method="MLE",
             response = res_cat, nugget=1e-12*var(res_cat), covtype = "exp")

gp_risk <- km(~.^2, design=runs[,c("Ftarget", "Btrigger")], estim.method="MLE",
              response = res_risk, nugget=1e-12*var(res_risk), covtype = "exp")
```

where `covtype` tells us the type of covariance function we are using, `estim.method` tells us how to estimate unknown parameters and `nugget` is an optional value accounting for the homogeneous nugget effect. The `design` variable gives the input parameters of the evaluated points we have so far (Roustant 2025). Here, `estim.method` is set to "MLE" which means we are using the maximum likelihood estimate to get the parameters μ_0 and cov_0 for our GPs.

We use a GP to emulate the objective function because it is much cheaper to evaluate compared to our objective function. If we let F and f be as in the Gaussian Processes section, then we can calculate $F(x)$ for any x in the design space as our estimate of $f(x)$ based on our current beliefs. This is true even for the evaluated points x_1, \dots, x_n as the emulator is fitted to these points(Spence 2025).

In the code, we predict the $\log(\text{catch})$ and $\log(\text{risk})$ at every point in the design space using the code below:

```
pred_risk_g <- predict(gp_risk,newdata=gridd,type="SK")
pred_cat_g <- predict(gp_cat,newdata=gridd,type="SK")
```

and then we can find the probability that the $\log(\text{risk})$ is ≤ 0.05 by using:

```
prisk <- pnorm(log(0.05),pred_risk_g$mean,pred_risk_g$sd+1e-12)
```

We can then exponentiate these results where needed. We are building our GPs with log of the values we want because this helps us generate better predictions(Huang et al. 2006).

We'll now explore how these GPs can be used to estimate the values of functions in the Kriging section.

Kriging

Kriging is a Bayesian statistical method for modelling functions(Frazier 2018). Let f be the objective function. We want to focus on the values $F(x_1), \dots, F(x_k)$ where $X := \{x_1, \dots, x_k\}$ is the finite design space. We can represent these values as the vector $F(x_{1:k}) := (F(x_1), \dots, F(x_k))$. This vector is unknown, but we can assume it was drawn at random from a multivariate normal distribution(Frazier 2018).

We let $\mu_0(x_{1:k})$ be the mean vector which we construct by evaluating μ_0 at each point in X . To encode the belief that as the inputs change slightly, the function only changes slightly, we choose cov_0 such that any points x_i, x_j that are close together in X have a large positive correlation(Frazier 2018). Then, we can model $F(x_{1:k})$ as below:

$$F(x_{1:k}) \sim N(\mu_0(x_{1:k}), \Sigma_{1:k}) \quad (2)$$

where

$$\mu_0(x_{1:k}) = (\mu_0(x_1), \dots, \mu_0(x_k)) \quad (3)$$

and N is the normal distribution(Frazier 2018). Now, if we have evaluated n points such that we have $f(x_{1:n})$ and want to evaluate x_{n+1} we let $k = n + 1$ in equation (2). Then, we can compute the conditional distribution of $F(x_{n+1})$ given $f(x_{1:n})$ using Bayes' rule:

$$F(x_{n+1})|f(x_{1:n}) \sim N(\mu_n(x_{n+1}), \sigma_n^2(x_{n+1})) \quad (4)$$

where:

$$\mu_n(x_{n+1}) = \text{cov}_0(x_{n+1}, x_{1:n})(\Sigma_{1:n})^{-1}(f(x_{1:n}) - \mu_0(x_{1:n})) + \mu_0(x_{n+1}) \quad (5)$$

$$\sigma_n^2(x_{n+1}) = \text{cov}_0(x_{n+1}, x_{n+1}) - \text{cov}_0(x_{n+1}, x_{1:n})(\Sigma_{1:n})^{-1} \text{cov}_0(x_{1:n}, x_{n+1}) \quad (6)$$

This conditional distribution $F(x_{n+1})|f(x_{1:n})$ is called the posterior probability distribution for x_{n+1} . We can calculate this distribution for every point in the design space X . This results in a new GP F_n with a mean vector and covariance kernel that depend on the location of the unevaluated points, the locations of the evaluated points $x_{1:n}$, and their values $f(x_{1:n})$ (Frazier 2018). So, we can use this in every round after updating it.

Our process for updating GPs can be seen in the code below:

```
gp_risk <- km(~.^2, design = runs[, c("Ftarget", "Btrigger")], estim.method = "MLE",
             response = res_risk, nugget = 1e-12 * var(res_risk), covtype = "exp")
```

This looks the same as before, but the `runs` variable includes all of our evaluated points and so we are adding the points that have been newly evaluated this round. hence, we can do the calculations above to update the GP with a new mean vector and covariance kernel. As we iterate over this, it will add the newly evaluated points for every round until `runs` becomes all of the points we have evaluated throughout the whole optimisation process.

Bayesian History Matching

Once the objective function f has been evaluated at eight points for the first round, Bayesian history matching speeds up the process by removing any points that are implausible(Spence 2025). This is done by using Bayes' Theorem to update our beliefs about the value of every unevaluated point based on the values of the points we have evaluated so far(Spence 2025). The general process is described below.

Let f be the objective function and x be a point in the sample space. We begin with some uncertainty about $f(x)$ (Spence 2025). However, we can make probabilistic statements such as:

$$P(f(x) > a) = \int_a^{\infty} P(f(x))df(x) \quad (7)$$

Once we evaluate another point x' where $x' \neq x$, we are able to improve our integral to:

$$P(f(x) > a|f(x')) = \int_a^{\infty} P(f(x)|f(x'))df(x) \quad (8)$$

We now let $a = \max\{f(x_1), \dots, f(x_l)\}$ where l is the number of points we have evaluated so far. For the first round, $l = 8$ but as the rounds increase we make sure to include all previous points of the objective function that have been evaluated.(Spence 2025) We remove the point x if:

$$P(f(x) > a|f(x_1), \dots, f(x_l)) = \int_a^{\infty} P(f(x)|f(x_1), \dots, f(x_l))df(x) < \varepsilon \quad (9)$$

using $\varepsilon = 0.00001$ until no plausible points remain. Then, the optimum will be x^* such that:

$$f(x^*) = \max\{f(x_1), \dots, f(x_l)\} \quad (10)$$

as our index l counts the number of points we have evaluated throughout the whole simulation.

In our case, x^* is the F_{target} and $B_{trigger}$ that will give the highest catch whilst following the precautionary principle(Spence 2025).

This process is written up in the code as the below:

```

#Finds the probability that log(risk) <= log(0.05)
prisk <- pnorm(log(0.05),pred_risk1_g$mean,pred_risk1_g$sd+1e-12)

#Finds the current maximum catch which also has risk < 0.05
#Note that we have exponentiated for this part
current_max <- max(runs$C_long[runs$risk3_long < 0.05])

#Finds the probability that the log(catch)<= log(current_max) for all points in the
#design space
pcat <- pnorm(log(current_max), pred_cat_g$mean, pred_cat_g$sd + 1e-12)

#Check which points have P(log(catch) >= log(curent_max)) > 0.00001 and have
#P(log(risk) <= log(0.05)) > 0.00001.
#These are our plausible points that remain.
possible <- (apply(cbind((1 - pcat), prisk), 1, min) > eps)

```

Expected Improvement

The first type of acquisition function we will look at is Expected Improvement (EI). Suppose we have sampled the points x_1, \dots, x_n and observe the values $f(x_{1:n})$. Then, if we were to return a solution at this point, bearing in mind we observe the objective function f without noise and we can only return points we have already evaluated, we would return $f_n^* = \max\{f(x_1), \dots, f(x_n)\}$. Imagine we then consider evaluating another point x_{n+1} to get $f(x_{n+1})$. As $f(x_{n+1})$ is unknown to us but we have our values $f(x_{1:n})$, we can say:

$$f(x_{n+1}) \sim F(x_{n+1})|f(x_{1:n})$$

We can then define the Expected Improvement as:

$$EI_n(x_{n+1}) := E_n[F(x_{n+1})|f(x_{1:n}) - f_n^*]^+ \quad (11)$$

where $[F(x_{n+1}) - f_n^*]^+$ is the positive part of $[F(x_{n+1}) - f_n^*]$ and E_n indicates the expectation taken under the posterior probability distribution defined in (4) so that we are using the previous points we have evaluated to help us find which new point is best to evaluate. This acquisition function is relatively easy to optimise and many different methods have been developed for doing this.

In our case, we want to sample multiple points at the same time, up to 8 points each round. This method of EI can be coded up as the function:

```

expected_improvement <- function(mu, sigma, y_best, xi = 0.05, task = "max",
                                  pred_risk, eps = 1e-4)
{
  #Setting up EI vector
  ei <- numeric(length(mu))
  #Determining points with P(risk<=0.05)>0.00001
  safe_points <- pred_risk >= eps

  # Only calculate EI if there are safe points and the task is a valid option
  if (any(safe_points)) {
    if (task == "min")
      imp <- y_best - mu[safe_points] - xi
    if (task == "max")
      imp <- mu[safe_points] - y_best - xi
    else
      stop('task must be "min" or "max"')
  }

  Z <- imp / sigma[safe_points]
  #Only calculate EI for safe points
  ei[safe_points] <- imp * pnorm(Z) + sigma[safe_points] * dnorm(Z)
  ei[safe_points][sigma[safe_points] == 0.0] <- 0.0
}

return(ei)
}

```

where `mu` is the mean of the GP for `log(catch)` and `sigma` is the standard deviation of the GP for `log(catch)`. Then we also have `y_best` which is the current best `log(catch)`, `xi` which is a parameter balancing exploration and exploitation. We set `pred_risk` = `prisk` and `eps` = `1e-4` and these are the same as before. This gives us the formula:

$$EI_n(x_{n+1}) = [(\mu_n(x_{n+1}) - f_n^* \cdot \Phi(Z)) + (\sigma_n(x_{n+1}) \cdot \phi(Z))]^+ \quad (12)$$

where again the notation $[.]^+$ means the positive part and where:

$$Z = \frac{\mu_n(x_{n+1}) - f_n^*}{\sigma_n(x_{n+1})}$$

Equation (12) can be gained from equation (2) by setting $k = n + 1$ and then studying the distribution of $f(x_{n+1}) - f_n^*$. However, we can also consider it as a version of equation (15) from (Jones, Schonlau, and Welch 1998) where we first flip the signs are we are focused on the maximisation case and then set $f_{min} = f_n^*$, $\hat{y} = \mu_n(x)$ and $s = \sigma_n(x)$.

Augmented Expected Improvement

By adjusting Equation (12) found in (Huang et al. 2006) to our own notation, we get that:

$$AEI_n(x_{n+1}) = EI_n(x_{n+1}) \left(1 - \frac{\sigma_{obs}}{\sqrt{\sigma_n^2(x_{n+1}) + \sigma_{obs}^2}} \right) \quad (13)$$

where σ_{obs} is the standard deviation of the noise variable set by the user and $\sigma_n(x)$ is the standard deviation of GP F at the n^{th} iteration, as used beforehand.

GUSTAV ASKS IF THSI IS RELATED TO THE METHOD OF ADDING THE VARIANCE TO THE COV MATRIX? In the current paper I'm looking at, it doesn't seem that they compare these. They get eqn (12) by integrating wrt $Y_p(x)$ over eqn(9). SEE IF I CAN FIND A SOURCE FOR THIS.

So, in the code we only slightly tweak the expected improvement function:

```
augmented_expected_improvement <- function(mu, sigma, y_best, xi = 0.05, task = "max",
                                             pred_risk, eps = 1e-4, noise_var = 0)
{
  ei <- numeric(length(mu))
  safe_points <- pred_risk >= eps

  # Only calculate AEI for safe points
  if (any(safe_points))
    if (task == "min")
      imp <- y_best - mu[safe_points] - xi
    if (task == "max")
      imp <- mu[safe_points] - y_best - xi
    else
      stop('task must be "min" or "max"')

  Z <- imp / sigma[safe_points]
  ei[safe_points] <- imp * pnorm(Z) + sigma[safe_points] * dnorm(Z)
  ei[safe_points][sigma[safe_points] == 0.0] <- 0.0

  # Augmentation factor to handle noise
  # When noise_var = 0, augmentation_factor = 1 (reduces to standard EI)
  augmentation_factor <- 1 - sqrt(noise_var / (noise_var + sigma^2))
  aei <- ei * augmentation_factor

  # Giving points with prob(risk <= 0.05) < 0.0001 an expected improvement
```

```

# of zero so we avoid them
aei <- ifelse(pred_risk < eps, 0, aei)

return(aei)
}

```

where `noise_var` corresponds to σ_{obs}^2 and `sigma` corresponds to $\sigma_n(\cdot)$.

Knowledge Gradient

We remove the assumption of EI that we have to return a pre-evaluated point as our best point(Frazier 2018). This allows us to do some different computations to the ones in EI. We also now start by saying that the solution we would choose if we have to stop sampling after n points would be the point in the design space with the largest $\mu_n(x)$ value, where μ_n is the mean vector of the posterior probability distribution after n iterations. We call this maximum x_n^* and then can say that $f(x_n^*)$ is random under the posterior distribution and has the conditional expected value:

$$\mu_n^* := \mu_n(x_n^*) = \max_x \mu_n(x) \quad (14)$$

where x is any point in the sample space(Frazier 2018).

Then, we imagine that we are now allowed to sample a new point x_{n+1} . We get a new posterior distribution which we can calculate using equation (4) by replacing x_{n+1} with x and $x_{1:n}$ with $x_{1:n+1}$ to include our new observation. This will have the posterior mean function $\mu_{n+1}(\cdot)$. The conditional expected value for $f(x_n^*)$ changes to be:

$$\mu_{n+1}^* := \max_x \mu_{n+1}(x) \quad (15)$$

So, we can see that the increase in the conditional expected value of $f(x_n^*)$ by sampling the new point x_{n+1} is(Frazier 2018):

$$\mu_{n+1}^* - \mu_n^* \quad (16)$$

While this quantity is unknown before we sample x_{n+1} we can calculate it's expected value given our observations x_1, \dots, x_n . The Knowledge Gradient for sampling at a new point x in the design space is defined as(Frazier 2018):

$$KG_n(x) := E_n[\mu_{n+1}^* - \mu_n^* | x_{n+1} = x] \quad (17)$$

where again E_n indicates the expectation taken under the posterior distribution at the n^{th} iteration. We would sample the point x with the largest $KG_n(x)$ as our next point[Frazier (2018)].

The easiest way to calculate the KG is via simulation. This can be done by simulating one possible value for $f(x_{n+1})$. Then, we calculate μ_{n+1}^* and subtract μ_n^* . We iterate this process many times so that we can find the average of $\mu_{n+1}^* - \mu_n^*$ and this allows us to estimate $KG_n(x)$ (Frazier 2018). This process, or calculating equation (17) directly from the properties of the normal distribution, both work well in discrete, low dimensional problems which is the situation we are in for the first half of the project(Frazier 2018).

Alternatively, we can calculate μ_{n+1} using the formula below(Ungredda, Pearce, and Branke 2022):

$$\mu_{n+1}(x) = \mu_n(x) + \frac{\text{cov}_n(x_{n+1}, x)}{\text{var}_n(x_{n+1}) + \sigma_{\text{obs}}^2}(F_{n+1} - \mu_n(x_{n+1})). \quad (18)$$

where F_{n+1} is the random distribution, a.k.a. the GP,

CHECK!

for $f(x_{n+1})$ which we tried to approximate above by simulating many times and σ_{obs}^2 is a noise variable which can be determined by the user(Ungredda, Pearce, and Branke 2022). From the GP for catch, we get the covariance matrix $\text{cov}_n(x_{n+1}, x)$ and the standard deviation σ_{obs} . Then, we have everything needed to compute equation (??).

This process is written as a function in the code:

```
knowledge_gradient_sim <- function(mu, sigma, model, obs_noise_var = 0, nsim = 100,
                                     pred_risk, eps = 1e-4)
{
  X_pred <- dat[, c("Ftrgt", "Btrigger")]
  cov_grid <- cov_exp(X_pred, X_pred, theta = model@covariance@range.val,
                        sigma2 = model@covariance@sd2)
  m <- length(mu)
  var <- sigma^2
  # Current best mean catch
  mu_best <- max(mu)
  kg <- numeric(m)

  # Loop over candidate points
  for (i in seq_len(m)) {
    # set knowledge gradient to 0 for any points with a predicted
    # risk that is too high
    if (pred_risk[i] < eps) {
```

```

    kg[i] <- 0
    next
}

mu_i <- mu[i]
sigma_i <- sqrt(var[i] + obs_noise_var)
# nsim simulated observations - using rnorm to sample from a
# normal distribution is ok here as we are running a GP
# to model the catch, which is assumed to be normally
# distributed
y_sim <- rnorm(nsim, mu_i, sigma_i)
# the covariance between every point in the grid and
# the point x_i
cov_xp_xi <- cov_grid[, i]
# denominator term in the KG update formula
denom <- var[i] + obs_noise_var
# For each simulated y, compute updated max(mu)
max_after <- numeric(nsim)
for (s in seq_len(nsim)) {
  # says what would the new maximum mu be if we observed this
  # simulated value at the point x_i
  # by implementing the standard formula for this in GP
  # posterior updating
  # This evaluates the posterior mean for every point in the
  # space, updating other points based on closeness
  # to the evaluated point by using the cov matrix
  mu_new <- mu + cov_xp_xi * (y_sim[s] - mu_i) / denom
  # takes this new maximum mu into a vector for later
  # averaging
  max_after[s] <- max(mu_new)
}
# Computes expected increase in the maximum posterior mean, mu
# higher values mean we're getting better knowledge as to where
# and what the maximum catch could be
kg[i] <- mean(max_after - mu_best)
}

# Returns the vector of knowledge gradient values for each point in
# the design space
kg
}

```

All the arguments of the function are as before, except that `obs_noise_var` is the new name

for the noise variable set by the user and `nsim` tells the code how many times to estimate the value of $f(x_{n+1})$ before we calculate $KG_n(x)$.

GUSTAV SAYS:

The main thing I would like you to add to the document is the choice of functional form of μ_0 and Σ_0 and the details of how you estimate the parameters in those functions. Are you doing that by maximum likelihood?

ADDED IN

If so, can you give the expression for the likelihood? Are you re-estimating the hyperparameters each time you evaluate the objective function at a new point?

THESE ARE BOTH THINGS I COULD DEFINETLY LOOK INTO.

References

- Frazier, Peter. 2018. “A Tutorial on Bayesian Optimization.” <https://doi.org/10.48550/arXiv.1807.02811>.
- Huang, D., Theodore Allen, William Notz, and Ning Zheng. 2006. “Global Optimization of Stochastic BlackBox Systems via Sequential Kriging Meta-Models.” *Journal of Global Optimization* 34: 441–66. <https://doi.org/10.1007/s10898-005-2454-3>.
- Jones, Donald, Matthias Schonlau, and William Welch. 1998. “Efficient Global Optimization of Expensive Black-Box Functions.” *Journal of Global Optimization* 13: 455–92. <https://doi.org/10.1023/A:1008306431147>.
- Roustant, Olivier. 2025. <https://www.rdocumentation.org/packages/DiceKriging/versions/1.6.1>.
- Spence, Michael A. 2025. “Using History Matching to Speed up Management Strategy Evaluation Grid Searches.” *Canadian Journal of Fisheries and Aquatic Sciences* 82: 1–14. <https://doi.org/10.1139/cjfas-2024-0191>.
- Ungredda, Juan, Michael Pearce, and Juergen Branke. 2022. “Efficient Computation of the Knowledge Gradient for Bayesian Optimization.” <https://arxiv.org/abs/2209.15367>.