

Supplementary Material

Deriving and optimising the cloaking variables

When solving for M , we put an arbitrary (positive) bound on Δ of 1, as any scaling of Δ caused by manipulating M , will scale the determinant of M by the inverse amount.⁶

We will express our problem as trying to *maximise* the entropy of a k -dimensional Gaussian distribution with covariance $P = M^{-1}$;

$$\begin{aligned} & \text{Maximise } \ln(|P|), \text{ subject to } n \text{ constraints,} \\ & 0 \leq \mathbf{c}_i^\top P \mathbf{c}_i \leq 1. \end{aligned}$$

Considering just the upper bounds and expressing this as a constraint satisfaction problem using Lagrange multipliers we have $L = \ln |P| + \sum_i \lambda_i (1 - \mathbf{c}_i^\top P \mathbf{c}_i)$. Differentiating (and setting to zero), $\frac{\partial L}{\partial P} = P^{-1} - \sum_i \lambda_i \mathbf{c}_i \mathbf{c}_i^\top = 0$. We also have the slackness conditions (for all i), $\lambda_i (\mathbf{c}_i^\top P \mathbf{c}_i - 1) = 0$ and we also require that $\lambda_i \geq 0$. Rearranging the derivative, we have, $P^{-1} = \sum_i \lambda_i \mathbf{c}_i \mathbf{c}_i^\top$. Note that as $\lambda_i \geq 0$, P^{-1} is positive semi-definite (psd),⁷ thus the initial lower bound (that $\mathbf{c}_i^\top P \mathbf{c}_i \geq 0$) is met. The upper bound (that $\mathbf{c}_i^\top P \mathbf{c}_i \leq 1$) is achieved if the λ_i are correctly chosen, such that the Lagrange and slackness conditions are met.

We now must maximise the expression for L wrt λ_j . To assist with this we rewrite our expression for $P^{-1} = \sum_i \lambda_i \mathbf{c}_i \mathbf{c}_i^\top = C \Lambda C^\top$, where $C = [\mathbf{c}_1, \mathbf{c}_2 \dots \mathbf{c}_n] = K_{*f} K^{-1}$ and Λ is a diagonal matrix of the values of λ_i . The summation in the expression for L , $\sum_i \lambda_i (1 - \mathbf{c}_i^\top P \mathbf{c}_i)$ can be written as $\text{Tr}(\Lambda - C^\top P C \Lambda)$. Substituting in our definition of P , we can write the summation as: $\text{Tr}(\Lambda - C^\top (C \Lambda C^\top)^{-1} C \Lambda) = \text{Tr}(\Lambda - C^\top C^{-\top} \Lambda^{-1} C^{-1} C \Lambda)$. Assuming C is invertible, the summation becomes $\text{Tr}(\Lambda - \Lambda^{-1} \Lambda)$. Differentiating this wrt λ_j equals one. We can use this result to find the

⁶For example if we write a new Δ' with M as mM , we see that we can take out the m term from Δ' 's inequality, leaving $\Delta' = m^{-1/2} \Delta$. When we scale Z , which has covariance mM by this new value of Δ' the covariance of the scaled Z equals $(\Delta')^2 mM = (m^{-1/2} \Delta)^2 mM = \Delta^2 M$, the magnitude change cancels, so any positive value of Δ (e.g. 1) will suffice for the optimisation. Also: in the following we ignore d and reintroduce it at the end of the derivation by defining $\Delta = d$ instead of it equalling 1.

⁷The summation is of a list of positive semi-definite rank-one matrices. One can see that such a sum is also positive semi-definite (to see this, consider distributing \mathbf{z}^\top and \mathbf{z} over the summation).

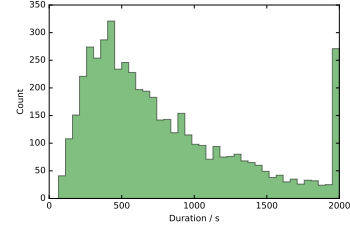


Figure 4: The duration of 5000 Citi Bike journeys. Note the effect caused by thresholding at 2000 s.

gradient of L wrt λ_j :

$$\begin{aligned} \frac{\partial L}{\partial \lambda_j} &= \frac{\partial \ln |P|}{\partial \lambda_j} + \frac{\partial}{\partial \lambda_j} \sum_i \lambda_i (1 - \mathbf{c}_i^\top P \mathbf{c}_i) \\ &= \text{Tr} \left(P^{-1} \frac{\partial P}{\partial \lambda_j} \right) + 1 \\ &= \text{Tr} (-P^{-1} P \mathbf{c}_j \mathbf{c}_j^\top P) + 1 \\ &= -\text{Tr} (\mathbf{c}_j \mathbf{c}_j^\top M^{-1}) + 1 \\ &= -\mathbf{c}_j^\top M^{-1} \mathbf{c}_j + 1 \end{aligned} \tag{10}$$

This can be solved using a gradient descent method, to give us those values of λ_i which minimise $\log |M|$ while ensuring $\Delta \leq 1$.

Citi Bike duration distribution

Figure 4 illustrates the distribution of journey times and the effect of thresholding.

Algorithms

Algorithm 2 describes the Cloaking method. Algorithm 1 describes the earlier method, from Section 2.

Hyperparameter selection

So far in this paper we have selected the values of the kernel hyperparameters *a priori*. Normally one may maximise the marginal likelihood to select the values or potentially integrates over the hyperparameters. In differential privacy we must take care when using private data to make these choices. Previous work exists to perform this selection, for example Kusner et al. [2015] describes a method for performing differentially private Bayesian optimisation, however their method assumes the training data is not private. Kusner et al. [2015] do suggest that the work of Chaudhuri and Vinterbo [2013] may allow Bayesian Optimisation to work in the situation in which the training data also needs to be private.

We decided instead that, due to the low-dimensionality of many hyperparameter problems, a simple grid search, combined with the exponential mechanism may allow the

Algorithm 1 Using the initial DP algorithm

Require: \mathcal{M} , the GP model (the kernel, hyperparameters and training inputs and normalised outputs)

Require: $X_* \in R^{P \times D}$, (the matrix of test inputs)

Require: $d > 0$, data sensitivity (maximum change possible) & $\varepsilon > 0, \delta \geq 0$, the DP parameters

```

1: function DIFFERENTIALLYPRIVATEREGRESSION( $X_*, M, d, \varepsilon, \delta$ )
2:    $K \leftarrow \mathcal{M}.\text{GET\_K}$  ▷ covariance between training points
3:    $\Delta \leftarrow d^2 \text{B}(K^{-1})^2$ 
4:    $\mathbf{f}_*, \sigma_*^2 \leftarrow \mathcal{M}.\text{GET\_PREDICTIONS}(X_*)$  ▷ Calculate non-DP predictions
5:    $\mathbf{f}_* \leftarrow \mathbf{f}_* + (\Delta d c(\delta) / \varepsilon) \mathbf{z}$  ▷  $\mathbf{z} \sim \mathcal{N}(0, K)$ 
6:   return  $\mathbf{f}_*, \sigma_*^2$ 
7: end function

8: function  $\text{B}(K^{-1})$ 
9:   return  $\max(|-[K_A^{-1}]^-|_\infty, |[K_A^{-1}]^+|_\infty)$ 
10: end function

11: function  $c(\delta)$ 
12:   return  $\sqrt{2 \log(2/\delta)}$ 
13: end function

```

selection of an acceptable set of hyperparameters. For the utility function we considered using the log marginal likelihood, with additional noise in the data-fit term to capture the effect of the DP noise. However for simplicity in estimating the bound and to avoid overfitting we simply used the sum square error (SSE) over a series of K -fold cross-validation runs, which for a given fold is $\sum_{i=1}^N (y_{*i} - y_{ti})^2$, with predictions \mathbf{y}_* and test values \mathbf{y}_t .

Before proceeding we need to compute a bound on the sensitivity of the SSE. To briefly recap, the DP assumption is that one data point has been perturbed by at most d . We need to bound the effect of this perturbation on the SSE. First we realise that this data point will only be in the *test set* in one of the K folds. In the remaining folds it will be in the training data.

If the perturbed data point is in the training data (\mathbf{y}), then we can compute the sensitivity of the SSE. The perturbation this would cause to the predictions (\mathbf{y}_*) is described using standard GP regression (and the cloaking matrix). Specifically a change of d in training point j will cause a $d\mathbf{c}_{jk}$ change in the test point predictions, where \mathbf{c}_{jk} is the j th column of the cloaking matrix for the k th fold.

To compute the perturbation caused by the change in the training data, we note that the SSE is effectively the square of the euclidean distance between the prediction and the test data. We are moving the prediction by $d\mathbf{c}_{jk}$. The largest effect that this movement of the prediction point could have on the distance between prediction and test locations is if it moves the prediction in the opposite direction to the test points. Thus it can increase (or decrease) the distance between the test and predictions by the largest length of $d\mathbf{c}_{jk}$ over training points. Hence

for one of the folds, the largest change in the SSE is $d^2 \max_j |\mathbf{c}_{jk}|_2^2$.

If the perturbed data point, j , was in the test data then the SSE will change by $(y_{*j} + d - y_{tj})^2 - (y_{*j} - y_{tj})^2 = d^2 + 2d(y_{*j} - y_{tj})$. The last part of the expression (the error in the prediction for point j) is unbounded. To allow us to constrain the sensitivity we enforce a completely arbitrary bound of being no larger than $\pm 4d$, thresholding the value if it exceeds this. Thus a bound on the effect of the perturbation is $d^2 + 2d \times 4d = d^2 + 8d^2 = 9d^2$.

The SSE of each fold is added together to give an overall SSE for the cross-validation exercise. We sum the $K - 1$ largest sensitivities and add $9d^2$ to account for the effect of the single fold in which the perturbing data point, j , will be in the test set. The perturbation could have been in the test data in any of the folds. We assumed it was in the fold with the smallest training-data sensitivity to allow us to make the result a lower bound on the sensitivity of the SSE to the perturbation. If it had been in any other fold the sensitivity would have been less (as more sensitivity would be contributed by the sum over the training sensitivities). Thus the sensitivity of the SSE over the K folds is; $9d^2 + \sum_{k=1}^{K-1} d^2 \max_j |\mathbf{c}_{jk}|_2^2$ (where the k folds are ordered by decreasing sensitivity)

We compute the SSE and the SSE's sensitivity for each of the hyperparameter combinations we want to test. We then use the computed sensitivity bound with the exponential mechanism to select the hyperparameters. To summarise, to use the exponential mechanism one evaluates the utility $u(x, r)$ for a given database x and for elements r , from a range. One also computes the sensitivity of this utility function by picking the highest sensitivity of any

Algorithm 2 Using the cloaking algorithm

Require: \mathcal{M} , the GP model (the kernel, hyperparameters and training inputs and normalised* outputs)

Require: $X_* \in R^{P \times D}$, (the matrix of test inputs)

Require: $d > 0$, data sensitivity (maximum change possible) & $\varepsilon > 0, \delta \geq 0$, the DP parameters

```

1: function DIFFERENTIALLYPRIVATECLOAKINGREGRESSION( $X_*, M, d, \varepsilon, \delta$ )
2:    $C \leftarrow \mathcal{M}.\text{GET\_C}(X_*)$  ▷ Compute the value of the cloaking matrix ( $K_{*f}K^{-1}$ )
3:    $\lambda \leftarrow \text{FINDLAMBDA}(C)$ 
4:    $M \leftarrow \text{CALCM}(\lambda, C)$  ▷ Calculate the DP noise covariance matrix
5:    $\Delta \leftarrow \text{CALCDELTA}(\lambda, C)^\dagger$ 
6:    $c \leftarrow \sqrt{2 \log \frac{2}{\delta}}$ 
7:    $y_*, \sigma_*^2 \leftarrow \mathcal{M}.\text{GET\_PREDICTIONS}(X_*)$  ▷ Calculate non-DP predictions
8:    $\tilde{y}_* \leftarrow y_* + (\Delta dc/\varepsilon)z$  ▷  $z \sim \mathcal{N}(0, M)$ 
9:   return  $\tilde{y}_*, \sigma_*^2$ 
10: end function

11: function  $\mathcal{M}.\text{GET\_C}(X_*)$ 
12:   From  $\mathcal{M}$  compute  $K_{*f}$  and  $K^{-1}$  ▷ Compute covariances between training and test points
13:    $C \leftarrow K_{*f}K^{-1}$ 
14:   return  $C$ 
15: end function

16: function FINDLAMBDA( $C$ )
17:    $\lambda \leftarrow \text{UNIFORM}(0.1, 0.9)$  ▷ Initialise randomly‡
18:    $\alpha \leftarrow 0.05$  ▷ Learning rate
19:   do
20:      $\frac{dL}{d\lambda} \leftarrow \text{CALCGRADIENT}(\lambda, C)$ 
21:      $\Delta_\lambda \leftarrow -\frac{dL}{d\lambda}\alpha$ 
22:      $\lambda \leftarrow \lambda + \Delta_\lambda$ 
23:   while  $\Delta_\lambda > 10^{-5}$ 
24:   return  $\lambda$ 
25: end function

26: function CALCGRADIENT( $\lambda, C$ )
27:    $M \leftarrow \text{CALCM}(\lambda, C)$ 
28:   for  $0 \leq j < N$  do ▷ N, number of columns in cloaking matrix, C.
29:      $[\frac{dL}{d\lambda}]_j \leftarrow -\text{Tr}(M^+ C_{:j} C_{:j}^\top) + 1$ 
30:   end for
31:   return  $\frac{dL}{d\lambda}$ 
32: end function

33: function CALCM( $\lambda, C$ )
34:    $M \leftarrow \sum_i \lambda_i C_{:i} C_{:i}^\top$ 
35:   return  $M$ 
36: end function

37: function CALCDELTA( $\lambda, C$ )
38:    $M \leftarrow \text{CALCM}(\lambda, C)$ 
39:    $\Delta \leftarrow \max_j C_{:j}^\top M^+ C_{:j}$ 
40:   return  $\Delta$ 
41: end function

```

*We assume the user will handle normalisation.

[†]Although we should have optimised M such that $\Delta \leq 1$, it may not have completely converged, so we compute the Δ bound for the value of M we have actually achieved.

[‡]We have found that occasionally the algorithm fails to converge. To confirm convergence we have found it useful to reinitialise and run the algorithm a few times.

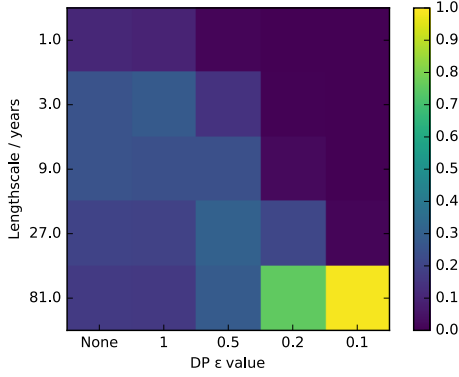


Figure 5: Effect of varying the differential privacy parameter, ϵ , on the likelihood of selecting each lengthscale. Colour indicates probability of parameter selection. With low privacy, a short lengthscale is appropriate which allows the GP to describe details in the data. With high privacy, a longer lengthscale is required, which will average over large numbers of individual data points.

of the utilities; in this case each utility corresponds to a negative SSE, and each sensitivity corresponds to the sum described above.

$$\Delta_u \triangleq \max_{r \in R} \max_{x, y} |u(x, r) - u(y, r)|$$

(where x and y are neighbouring databases).

The exponential mechanism selects an element r with probability proportional to:

$$\exp\left(\frac{\epsilon u(x, r)}{2\Delta_u}\right)$$

Note that for a given privacy budget, some will need to be expended on this selection problem, and the rest expended on the actual regression.

Effect of privacy on optimum hyperparameter selection

A final interesting result is in the effect of the level of privacy in the regression stage on the selection of the lengthscale. This is demonstrated in the distribution of probabilities over the lengthscales when we adjust ϵ . Figure 5 demonstrates this effect. Each column is for a different level of privacy (from none to high) and each tile shows the probability of selecting that lengthscale. For low privacy, short lengthscales are acceptable, but as the privacy increases, averaging over more data allows us to give more accurate answers.

Privacy on the training inputs

To release the mean function such that the training inputs remain private, we need a general bound on the infinity norm of the covariance function, that does not depend explicitly on the values of X .

Varah [1975] show that if J is *strictly diagonally dominant*⁸ then:

$$\|J^{-1}\|_{\infty} \leq \max_{1 \leq i \leq n} \frac{1}{\Delta_i(J)} = b(J)$$

where we have defined this bound as $b(J^{-1})$. We also define $\Delta_i(J) = |J_{ii}| - \sum_{j \neq i} |J_{ij}|$, i.e. the sum of the off diagonal elements in row i subtracted from the diagonal element.

So if K is strictly diagonally dominant (which is achieved if the inputs are sufficiently far apart, and/or if sufficient uncertainty exists on the outputs), then we have a bound on the sums of its rows. The above bound means that,

$$\sum_{i=1}^n \alpha_i - \alpha'_i \leq \Delta_y b(J^{-1}) \quad (11)$$

To ensure sufficient distance between inputs, we could use inducing variables, which can be arbitrarily placed, so that the above constraints on the covariance matrix are observed.

Integral Kernel

We used a custom kernel in the Citi Bike comparison to make predictions over the binned DP-noisy dataset. In summary the observations were considered to be the integrals over each bin (provided by multiplying the noisy means by the sizes of the bins). The predictions were made at points on the latent function being integrated. This allowed us to make predictions from the binned data in a principled manner.

Considering just one dimension. To compute the covariance for the integrated function, $F(\cdot)$, we integrate the original EQ kernel, $k_{ff}(\cdot, \cdot)$, over both its input values,

$$k_{FF}((s, t), (s', t')) = \alpha \int_s^t \int_{s'}^{t'} k_{ff}(u, u') du' du,$$

so that given two pairs of input locations, corresponding to two definite integrals, we can compute the covariance between the two.

Similarly we can calculate the cross-covariance k_{Ff} between F and f . Both k_{FF} and k_{Ff} can be extended to multiple dimensions. Each kernel function contains a unique lengthscale parameter, with the bracketed kernel subscript

⁸A matrix, J , is strictly diagonally dominant if $\Delta_i(J) > 0$ for all $1 \leq i \leq n$.

index indicating these differences. We can express the new kernel as the product of our one dimensional kernels:

$$k_{FF} = \prod_i k_{FF(i)}((s_i, t_i), (s'_i, t'_i)),$$

with the cross covariance given by

$$k_{Ff} = \prod_i k_{Ff(i)}((s_i, t_i), (s'_i, t'_i)).$$

The above expressions can then be used to make predictions of the latent function f given observations of its definite integrals.