
Local Gaussian Regression

Franziska Meier¹Philipp Hennig²Stefan Schaal^{1,2}¹Computational Learning and Motor Control Lab, University of Southern California, Los Angeles²Max-Planck-Institute for Intelligent Systems, Tübingen, Germany

Abstract

Locally weighted regression was created as a nonparametric learning method that is computationally efficient, can learn from very large amounts of data and add data incrementally. An interesting feature of locally weighted regression is that it can work with spatially varying length scales, a beneficial property, for instance, in control problems. However, it does not provide a generative model for function values and requires training and test data to be generated identically, independently. Gaussian (process) regression, on the other hand, provides a fully generative model without significant formal requirements on the distribution of training data, but has much higher computational cost and usually works with one global scale per input dimension. Using a localising function basis and approximate inference techniques, we take Gaussian (process) regression to increasingly localised properties and toward the same computational complexity class as locally weighted regression.

1 Introduction

Besides expressivity and sample efficiency, computational cost is a crucial design criterion for machine learning algorithms in real-time settings, such as control problems. An example is the problem of building a model for robot dynamics: The sensors in a robot's limbs can produce thousands of datapoints per second, quickly amassing a local coverage of the input domain. In such settings, fast local learning and generalization can be more important than a globally optimized model. A learning method should rapidly produce a good model from the large number N of datapoints, using a comparably small number M of parameters.

Locally weighted regression (LWR) [1] makes use of the popular and well-studied idea of local learning [2] to address the task of compressing large amounts of data into a small number of parameters. In the spirit of a Taylor expansion, the idea of LWR is that simple models with few parameters may locally be precise,

while it may be difficult to find good nonlinear features to capture the entire function globally – lots of good local models may form a good global one.

The key to LWRs low computational cost (linear, $\mathcal{O}(NM)$) is that each local model is trained independently. The resulting speed has made LWR popular in robot learning. The downside is that LWR requires several tuning parameters, whose optimal values can be highly data dependent. This is at least partly a result of the strongly localized training strategy, which does not allow models to ‘coordinate’, or to benefit from other local models in their vicinity.

Here, we explore a probabilistic alternative to LWR that alleviates the need for parameter tuning, but retains potential for fast training. An initial candidate could be the mixture of experts model (ME) [3]. Indeed, it has been argued [1], that LWR can be thought of as a mixture of experts model in which experts are trained independently of each other. The advantage of ME is that it comes with a full generative model [4, 5] allowing for principled learning of parameters and expert allocations, while LWR does not (see Figure 1). However, in this work we emulate LWRs assumption that the data is a continuous function as opposed to a mixture of linear models, for instance. When the underlying data is not assumed to be generated by a mixture model, utilizing ME to fit the data makes inference (unnecessarily) complicated. Hence, we are interested in a probabilistic model that captures the best of both worlds without making the mixture assumption: A generative model that has the ability to localize computations to speed up learning.

Our proposed solution is *local Gaussian regression* (LGR), a linear regression model that explicitly encodes localisation. Starting from the well-known probabilistic, Gaussian formulation of least-squares generalized linear regression, we first re-interpret the well known radial basis feature functions as localisers of constant feature functions. This allows us to introduce more expressive features, capturing the idea of local models within the Gaussian regression framework. In its exact form, this model has the cubic cost in M typical of

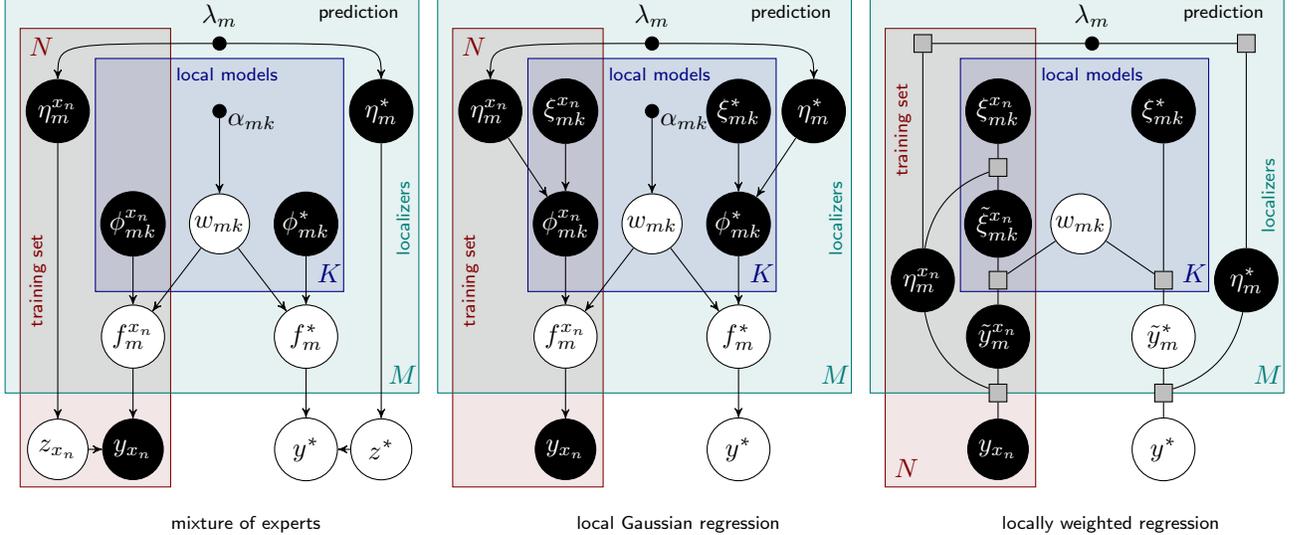


Figure 1: Generative model for mixture of experts (ME) (left), local Gaussian regression (LGR) (middle), and factor graph for locally weighted regression (LWR) (right). LWR assumes a fixed contribution of each local model to each observation. So a V-structure only exists *within* each m -plate, making inference cheaper. In ME, contributions of experts to each observations are normalized, coupling all the experts and making it necessary to update responsibilities in the E-step in a global manner. In local Gaussian regression, the contribution of each model to observations is uncertain (similar to ME), so the model is more densely connected than LWR. However, LWR is not a generative model for the data: It treats training and test data in different ways.

Gaussian models, arising because observations induce correlation between all local models in the posterior. To decouple the local models, we propose a variational approximation that gives essentially linear cost in the number of models M . The core of this work revolves around fitting the model parameters using maximum likelihood (Section 3.1). As a final note, we show how the novel feature function representation allows us to readily extend our probabilistic model to a local non-parametric formulation (Section 4).

Previous work on probabilistic formulations of local regression [6, 7] has been focussed on bottom-up constructions, trying to find generative models for one local model at a time. To our knowledge, this is the first top-down approach, starting from a globally optimal training procedure, to find approximations giving a localized regression algorithm similar in spirit to LWR.

2 Background

Both LWR and Gaussian regression have been studied extensively before, so we only give brief introductions here. Generalized linear regression maps weights $\mathbf{w} \in \mathbb{R}^F$ to the nonlinear function $f: \mathbb{R}^D \rightarrow \mathbb{R}$ via F feature functions $\phi_i(x): \mathbb{R}^D \rightarrow \mathbb{R}$:

$$f(x) = \sum_{i=1}^F \phi_i(x) w_i = \boldsymbol{\phi}^\top \mathbf{w}. \quad (1)$$

Using the feature matrix $\Phi \in \mathbb{R}^{N \times F}$ whose elements are $\Phi_{nf} = \phi_f(x_n)$, the function values at N locations $\mathbf{x}_n \in \mathbb{R}^D$, subsumed in the matrix $X \in \mathbb{R}^{N \times D}$, are $\mathbf{f}(X) = \Phi \mathbf{w}$.

Locally weighted regression (LWR) trains M local regression models. We will assume each local model has K local feature functions $\xi_{mk}(x)$, so that the m -th model's prediction at x is

$$f_m(x) = \sum_{k=1}^K \xi_{mk}(x) w_{mk} = \boldsymbol{\xi}_m(x) \mathbf{w}_m \quad (2)$$

$K=2$ and $\xi_{m1}(x) = 1, \xi_{m2}(x) = (x - c_m)$ gives a linear model around c_m . The models are localized by a non-negative, symmetric and integrable weighting $\eta_m(x)$, typically the radial basis function

$$\eta_m(x) = \exp \left[-\frac{(x - c_m)^2}{2\lambda_m^2} \right], \quad \text{or, for } \mathbf{x} \in \mathbb{R}^D, \quad (3)$$

$$\eta_m(\mathbf{x}) = \exp \left[-\frac{1}{2} (\mathbf{x} - \mathbf{c}_m) \Lambda^{-1} (\mathbf{x} - \mathbf{c}_m)^\top \right]. \quad (4)$$

with center c_m and length scale λ or positive definite metric Λ . In LWR, each local model is trained independently of the others, by minimizing a quadratic

loss

$$\begin{aligned} \mathbb{L}(\mathbf{w}) &= \sum_{n=1}^N \sum_{m=1}^M \eta_m(x_n) (y_n - \xi_m(x_n) \mathbf{w}_m)^2 \quad (5) \\ &= \sum_{n=1}^N \sum_{m=1}^M \left(\sqrt{\eta_m(x_n)} y_n - \sqrt{\eta_m(x_n)} \xi_m(x_n) \mathbf{w}_m \right)^2 \end{aligned}$$

over the N observations (y_n, x_n) . At test time, local predictions (2) are combined into a joint prediction at input x as a normalised weighted sum

$$f(x) = \frac{\sum_m \eta_m(x) f_m(x)}{\sum_{m'} \eta_{m'}(x)} = \sum_{m,k} \phi_{mk}(x) w_{mk} \quad (6)$$

with the features

$$\phi_{mk}(x) = \frac{\eta_m(x) \xi_{mk}(x)}{\sum_{m'} \eta_{m'}(x)}. \quad (7)$$

An important observation is that the objective (5) cannot be interpreted as least-squares estimation of the linear model Φ from Eq. (6). Eq. (5) is effectively training M linear models with $\phi_{mk}(x_n) = \sqrt{\eta_m(x_n)} \xi_{mk}(x_n)$ on M *separate* datasets $y_m(x_n) = \sqrt{\eta_m(x_n)} y_n$, but that model differs from the one in Equation (6). Thus, LWR can not be cast probabilistically as one generative model for training and test data simultaneously. The factor graph in Figure 1, right, shows this broken symmetry.

Training LWR is linear in N and M , and cubic in K , since it involves a least-squares problem in the K weights \mathbf{w}_m . But this low cost profile also means absence of mixing terms between the M models, thus no ‘coordination’ between the local models. One problem this causes is that model-specific hyperparameters are not sufficiently identified and tend to over-fit. To pick one example among the hyperparameters: when learning the length scale parameters λ_m , there is no sense of global improvement. Each local model focusses only on fitting data points within a λ_m -neighborhood around c_m . So the optimal choice is actually $\lambda \rightarrow 0$, which gives a useless regression model. To counteract this behaviour, implementations of LWR use regularisation of λ , but of course this introduces free parameters, whose best values are not obvious. Most works in LWR address this problem via leave-one-out cross-validation, which works well if the data is sparse. But when the data is dense, like in robotics, the kernel tends to shrink to only ‘activate’ points close by the left-out data point.

Gaussian regression [8, §2] is a probabilistic framework for inference on the weights \mathbf{w} given the features ϕ . Using a Gaussian prior $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}_0, \Sigma_0)$ and a Gaussian likelihood $p(\mathbf{y} | \phi, \mathbf{w}) = \mathcal{N}(\mathbf{y}; \phi^\top \mathbf{w}, \beta^{-1} \mathbf{I})$,

the posterior distribution is itself Gaussian:

$$p(\mathbf{w} | \mathbf{y}, \phi) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}_N, \Sigma_N) \quad \text{with} \quad (8)$$

$$\boldsymbol{\mu}_N = (\Sigma_0^{-1} + \beta \Phi^\top \Phi)^{-1} (\beta \Phi^\top \mathbf{y} - \Sigma_0^{-1} \boldsymbol{\mu}_0) \quad (9)$$

$$\Sigma_N = (\Sigma_0^{-1} + \beta \Phi^\top \Phi)^{-1} \quad (10)$$

The mean of this Gaussian posterior is identical with the Σ -regularised least-squares estimator for \mathbf{w} , so it could alternatively be derived as the minimiser of the loss function

$$\mathbb{L}(\mathbf{w}) = \mathbf{w}^\top \Sigma^{-1} \mathbf{w} + \sum_{n=1}^N (y_n - \phi(x_n) \mathbf{w})^2. \quad (11)$$

But the probabilistic interpretation of Equation (8) has additional value over (11) because it is a generative model for all (training and test) data points \mathbf{y} , which can be used to learn hyperparameters of the feature functions. The prediction for $f(x_*)$ with features $\phi(x_*) =: \phi_*$ is also Gaussian, with

$$p(f(x_*) | \mathbf{y}, \phi) = \mathcal{N}(f(x_*); \phi_* \boldsymbol{\mu}_N, \phi_* \Sigma_N \phi_*^\top) \quad (12)$$

As is widely known, this framework can be extended to the nonparametric case by a limit which replaces all inner products $\phi(x_1) \Sigma_0 \phi(x_2)^\top$ with a Mercer (positive semi-definite) kernel function $k(x_1, x_2)$. The corresponding priors are Gaussian processes. This generalisation will only play a role in Section 4 of this paper, but the direct connection between Gaussian regression and the elegant theory of Gaussian processes is often cited in favour of this framework. Its main downside, relative to LWR, is computational cost: Calculating the posterior (12) requires solving the least-squares problem for all F parameters \mathbf{w} jointly, by inverting the Gram matrix $(\Sigma_0^{-1} + \beta \Phi^\top \Phi)$. In the general case, this inversion requires $\mathcal{O}(F^3)$ operations. The point of the construction in the following sections is to use approximations to lower the computation cost of this operation such that the resulting algorithm is comparable in cost to LWR, while retaining the probabilistic interpretation, and the modelling robustness of the full Gaussian model.

The Gaussian regression framework puts virtually no limit on the choice of basis functions ϕ , even allowing discontinuous and unbounded functions, but the radial basis function (RBF, aka. Gaussian, square-exponential) features from Equation (3), $\phi_i(x) = \eta_i(x)$, (for $F = M$) enjoy particular popularity for various reasons, including algebraic conveniences and the fact that their associated reproducing kernel Hilbert space lies dense in the space of continuous functions [9]. A downside of this covariance function is that it uses one global length scale for the entire input domain. There are some special kernels of locally varying regularity [10], and mixture descriptions offer a more discretely varying model class [11]. Both, however, require computationally demanding training.

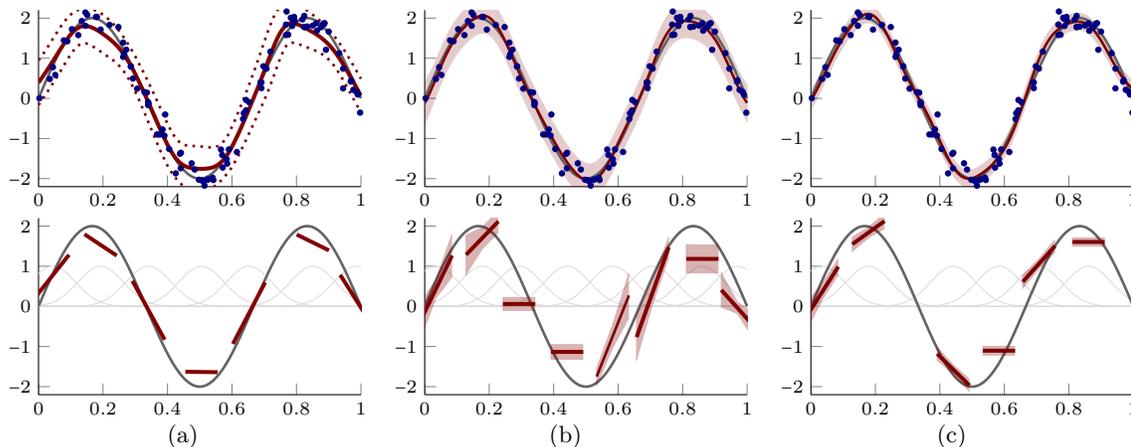


Figure 2: Noisy data drawn from sine function learned by (a) LWR (b) exact and (c) approximate local Gaussian regression. In LWR, local models do not know of each other and thus fit tangential lines around the center of each local model. In local Gaussian regression local models are correlated. They work together to fit the function. In the approximate version (c) this correlation is reduced, but not completely gone.

3 Local Parametric Gaussian Regression

In Gaussian regression with RBF features as described above, without changing anything, the features $\phi_m(x) = \eta_m(x)$ can be interpreted as M constant function models $\xi_m(x) = 1$, localised by the RBF function, $\phi(x) = [\xi_1(x)\eta_1(x), \dots, \xi_M(x)\eta_M(x)]$. This representation extends to more elaborate local models $\xi(x)$. For example $\xi(x) = x$ gives a local weighted regression with linear local models. Extending to M local models consisting of K parameters each, feature function ϕ_{mk}^n combines the k^{th} component of the local model $\xi_{km}(x_n)$, localised by the m -th weighting function $\eta_m(x)$

$$\phi_{mk}^n := \phi_{mk}(x_n) = \eta_m(x_n)\xi_{km}(x_n). \quad (13)$$

For these features, treating mk as the index set of a vector with MK elements, the results from Equations (8) and (12) apply, giving a localised linear Gaussian regression algorithm.

The choice of local parametric model is essentially arbitrary. But regular features are an obvious choice. A scalar polynomial model of order $K - 1$ arises from $\xi_{km}(x_n) = (x_n - c_m)^{k-1}$. Local linear regression in a K -dimensional input space takes the form $\xi_{km}(x_n) = x_n^k - c_m^k$. We will adopt this latter model in the remainder, because it is also the most widely used model in LWR. An illustration of LWR regression and the proposed local Gaussian regression is given in Figure 2.

Since it will become necessary to prune out unnecessary parts of the model, we adopt the classic idea of automatic relevance determination [12, 13] using a

factorizing prior

$$p(\mathbf{w}|A) = \prod_{m=1}^M \mathcal{N}(\mathbf{w}_m; 0, A_m^{-1}) \quad \text{with} \quad (14)$$

$$A_m = \text{diag}(\alpha_{m1}, \dots, \alpha_{mK}). \quad (15)$$

So every component k of each local model m has its own precision, and can thus be pruned out by setting $\alpha_{mk} \rightarrow \infty$. Section 3.1 assumes a fixed number M of local models with fixed centers c_m . Thus the parameters are $\theta = \{\beta, \{\alpha_{mk}\}, \{\lambda_{mk}\}\}$. We propose an approximation for estimating these parameters. Section 3.2 then describes placing the local models incrementally to adapt M and c_m .

3.1 Learning in Local Bayesian Linear Regression

Exact inference in Gaussian regression with localised feature functions comes with the cubic cost of its non-local origins. However, because of the localised feature functions, correlation between far away local models is approximately 0, hence inference is approximately independent between local models. In this section we aim to make use of this ‘‘almost independence’’ to derive a localised approximate inference scheme of low cost, similar in spirit to LWR. To arrive at this localised learning algorithm we first introduce a latent variable f_{nm} for each local model m and data point x_n , similar to probabilistic backfitting [14]. Intuitively, the \mathbf{f} form approximate targets, one for each local model, against which the local regression parameters fit (see also Figure 1, middle).

This modified model motivates a factorizing variational bound constructed in Section 3.1.1, rendering the the

local models computationally independent, which allows for fast approximate inference in the local Gaussian model. Hyperparameters will be learnt by approximate maximum likelihood (3.1.2), i.e. iterating between constructing a bound $q(\mathbf{z}|\boldsymbol{\theta})$ on the posterior over variables \mathbf{z} given current parameter estimates $\boldsymbol{\theta}$ and optimising q with respect to $\boldsymbol{\theta}$.

3.1.1 Variational Bound

The complete data likelihood of the modified model is

$$p(\mathbf{y}, \mathbf{f}, \mathbf{w} | \Phi, \theta) = \prod_{n=1}^N p(y_n | \mathbf{f}_n, \beta_y) \quad (16)$$

$$\prod_{n=1}^N \prod_{m=1}^M p(f_{nm} | \phi_m^n \mathbf{w}_m, \beta_{fm}) \prod_{m=1}^M p(\mathbf{w}_m | A_m)$$

(c.f. Figure 1, left). Our Gaussian model involves the latent variables \mathbf{w} and \mathbf{f} , the precision β and the model parameters λ_m, c_m . We treat \mathbf{w} and \mathbf{f} as probabilistic variables and estimate $\theta = \{\beta, \boldsymbol{\lambda}, \mathbf{c}\}$. On \mathbf{w}, \mathbf{f} , we construct a variational bound $q(\mathbf{w}, \mathbf{f})$ imposing factorisation $q(\mathbf{w}, \mathbf{f}) = q(\mathbf{w})q(\mathbf{f})$. The variational free energy is a lower bound on the log evidence for the observations y :

$$\log p(\mathbf{y} | \theta) \geq \int q(\mathbf{w}, \mathbf{f}) \log \frac{p(\mathbf{y}, \mathbf{w}, \mathbf{f} | \theta)}{q(\mathbf{w}, \mathbf{f})}. \quad (17)$$

This bound is maximized by the $q(\mathbf{w}, \mathbf{f})$ minimizing the relative entropy $D_{\text{KL}}[q(\mathbf{w}, \mathbf{f}) \| p(\mathbf{w}, \mathbf{f} | \mathbf{y}, \theta)]$, the distribution for which $\log q(\mathbf{w}) = \mathbb{E}_{\mathbf{f}}[\log p(\mathbf{y} | \mathbf{f}, \mathbf{w}) p(\mathbf{w}, \mathbf{f})]$ and $\log q(\mathbf{f}) = \mathbb{E}_{\mathbf{w}}[\log p(\mathbf{y} | \mathbf{f}, \mathbf{w}) p(\mathbf{w}, \mathbf{f})]$. It is relatively straightforward to show (e.g. [15]) that these distributions are Gaussian in both \mathbf{w} and \mathbf{f} . The approximation on \mathbf{w} is

$$\log q(\mathbf{w}) = \mathbb{E}_{\mathbf{f}} [\log p(\mathbf{f}_n | \phi(x_n), \mathbf{w}) + \log p(\mathbf{w} | A)]$$

$$= \log \prod_{m=1}^M \mathcal{N}(\mathbf{w}_m; \boldsymbol{\mu}_{w_m}, \Sigma_{w_m}) \quad (18)$$

where

$$\Sigma_{w_m} = \left(\beta_{fm} \sum_{n=1}^N \phi_m^n \phi_m^{nT} + A_m \right)^{-1} \in \mathbb{R}^{K \times K} \quad (19)$$

$$\boldsymbol{\mu}_{w_m} = \beta_{fm} \Sigma_{w_m} \left(\sum_{n=1}^N \phi_m^n \mathbb{E}[f_{nm}] \right) \in \mathbb{R}^{K \times 1} \quad (20)$$

The posterior update equations for the weights are local: each of the local models updates its parameters independently. This comes at the cost of having to update the belief over the variables f_{nm} , which achieves a coupling between the local models. The Gaussian variational bound on \mathbf{f} is

$$\log q(\mathbf{f}_n) = \mathbb{E}_{\mathbf{w}} [\log p(y_n | \mathbf{f}_n, \beta_y) + \log p(\mathbf{f}_n | \phi_m^n, \mathbf{w})]$$

$$= \mathcal{N}(\mathbf{f}_n; \boldsymbol{\mu}_{f_n}, \Sigma_f) \quad (21)$$

where

$$\Sigma_f = B^{-1} - B^{-1} \mathbf{1} (\beta_y^{-1} + \mathbf{1}^T B^{-1} \mathbf{1})^{-1} \mathbf{1}^T B^{-1}$$

$$= B^{-1} - \frac{B^{-1} \mathbf{1} \mathbf{1}^T B^{-1}}{\beta_y^{-1} + \mathbf{1}^T B^{-1} \mathbf{1}} \quad (22)$$

$$\boldsymbol{\mu}_{f_n} = \sum_{m=1}^M \mathbb{E}_{\mathbf{w}} [\mathbf{w}_m^T] \phi(x_n) +$$

$$\frac{1}{\beta_y^{-1} + \mathbf{1}^T B^{-1} \mathbf{1}} B^{-1} \mathbf{1} \left(y_n - \sum_{m=1}^M \mathbb{E}[\mathbf{w}_m]^T \phi_m^n \right) \quad (23)$$

where $\boldsymbol{\mu}_{f_n} \in \mathbb{R}^M$ and using $B = \text{diag}(\beta_{f1}, \dots, \beta_{fM})$. These updates can be performed in $\mathcal{O}(MK)$. Inspecting Equations (23) and (20), one can see that the optimal assignment of all $\boldsymbol{\mu}_w$ actually amounts to a joint linear problem of size $MK \times MK$, which could also be solved as a linear program. However, this would come at additional computational cost.

3.1.2 Optimizing Hyperparameters

We set the model parameters $\theta = \{\beta_y, \{\beta_{fm}, \lambda_m\}_{m=1}^M, \{\alpha_{mk}\}\}$ to maximize the expected complete log likelihood under the variational bound,

$$\mathbb{E}_{\mathbf{f}, \mathbf{w}} [\log p(\mathbf{y}, \mathbf{f}, \mathbf{w} | \Phi, \theta)] = \quad (24)$$

$$\mathbb{E}_{\mathbf{f}, \mathbf{w}} \left\{ \sum_{n=1}^N \left[\log \mathcal{N} \left(y_n; \sum_{m=1}^M f_{nm}, \beta_y^{-1} \right) \right. \right.$$

$$\left. \left. + \sum_{m=1}^M \log \mathcal{N} (f_{nm}; \mathbf{w}_m^T \phi_m^n, \beta_{fm}^{-1}) \right] \right.$$

$$\left. \left. + \sum_{m=1}^M \log \mathcal{N} (\mathbf{w}_m; 0, A_m^{-1}) \right\}$$

Setting the gradient of this expression to zero leads to the following update equations for the variances

$$\beta_y^{-1} = \frac{1}{N} \sum_{n=1}^N (y_n - \mathbf{1} \boldsymbol{\mu}_{f_n})^2 + \mathbf{1}^T \Sigma_f \mathbf{1} \quad (25)$$

$$\beta_{fm}^{-1} = \frac{1}{N} \sum_{n=1}^N \left[(\mu_{f_{nm}} - \boldsymbol{\mu}_{w_m}^T \phi_m^n)^2 + \phi_m^{nT} \Sigma_{w_m} \phi_m^n \right] + \sigma_{fm}^2 \quad (26)$$

$$\alpha_{mk}^{-1} = \mu_{w_{mk}}^2 + \Sigma_{w, kk} \quad (27)$$

The gradient with respect to the scales of each local model is completely localized

$$\frac{\partial \mathbb{E}_{\mathbf{f}, \mathbf{w}} [\log p(\mathbf{y}, \mathbf{f}, \mathbf{w} | \Phi, \theta)]}{\partial \log \lambda_{mk}}$$

$$= \frac{\partial \mathbb{E}_{\mathbf{f}, \mathbf{w}} \left[\sum_{n=1}^N \mathcal{N} (f_{nm}; \mathbf{w}_m^T \phi_m(x_n), \beta_{fm}^{-1}) \right]}{\partial \lambda_{mk}} \quad (28)$$

We use gradient ascent to optimize the length scales λ_{mk} . All necessary equations are of low cost and, with

the exception of the variance $\frac{1}{\beta_y}$, all hyper-parameter updates are solved independently for each local model, similar to LWR. In contrast to LWR, however, these local updates do not cause a shrinking in the length scales: In LWR, both inputs and outputs are weighted by the localizing function, thus reducing the length scale improves the fit. The localization of Equation (13) only affects the influence of regression model m , but the targets still need to be fit accordingly. Shrinking of local models only happens if it actually improves the fit against the unweighted targets f_{nm} .

3.1.3 Prediction

Prediction at a test point x_* arise from marginalizing over both \mathbf{w} and \mathbf{f} , using

$$\begin{aligned} & \int \mathcal{N}(y_*; \mathbf{1}^T \mathbf{f}_*, \beta_y^{-1}) \mathcal{N}(\mathbf{f}_*; W^T \phi(x_*), B^{-1}) d\mathbf{f}_* \\ &= \mathcal{N}(y_*; \sum_m \mathbf{w}_m^T \phi_m^*, \beta_y^{-1} + \mathbf{1}^T B^{-1} \mathbf{1}) \end{aligned} \quad (29)$$

and

$$\begin{aligned} & \int \mathcal{N}(y_*; \mathbf{w}^T \phi^*, \beta_y^{-1} + \mathbf{1}^T B^{-1} \mathbf{1}) \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}_w, \Sigma_w) d\mathbf{w} \\ &= \mathcal{N}\left(y_*; \sum_m \boldsymbol{\mu}_{w_m}^T \phi_m^*, \sigma^2(x^*)\right) \end{aligned} \quad (30)$$

where

$$\sigma^2(x^*) = \beta_y^{-1} + \sum_{m=1}^M \beta_{f_m}^{-1} + \sum_{m=1}^M \phi_m^{*T} \Sigma_{w_m} \phi_m^* \quad (31)$$

which is linear in M and K .

3.2 Incremental Adding of Local Models

Up to here, the number M and locations \mathbf{c} of local models were fixed. An extension analogous to the incremental learning of the relevance vector machine [16] can be used to iteratively add local models at new, greedily selected locations c_{M+1} . The resulting algorithm starts with one local model, and per iteration adds one local model in the variational step. Conversely, existing local models for which all components $\alpha_{mk} \rightarrow \infty$ are pruned out. This works well in practice, with the caveat that the number local models M can grow fast initially before the pruning becomes effective. Thus we check for each selected location c_{M+1} whether any of the existing local models $c_{1:M}$ produces a localizing weight $\eta_m(c_{M+1}) \geq w_{\text{gen}}$, where w_{gen} is a parameter between 0 and 1 and regulates how many parameters are added. An overview of the final algorithm is given in 1.

Algorithm 1 Incremental LGR

Require: $\{x_n, y_n\}_{n=1}^N$
 // initialize first local model
 1: $M = 1; c_1 = x_1; C = \{c_1\}$
 // iterate over all data points
 2: **for** $n = 2 \dots N$ **do**
 3: **if** $\eta_m(x_n) < w_{\text{gen}}, \forall m = 1, \dots, M$ **then**
 4: $c_m \leftarrow x_n$
 5: $C \leftarrow C \cup \{c_m\}, M = M + 1$
 6: **end if**
 // E-Step: equations (19),(20),(22),(23)
 7: $\{\boldsymbol{\mu}_{w_m}, \Sigma_{w_m}, \boldsymbol{\mu}_{f_m}, \Sigma_{f_m}\}_m \leftarrow$
 8: E-STEP($\{\beta_y, \beta_{f_m}, \{\alpha_{mk}, \lambda_{mk}\}_k\}_m$)
 // M-Step: equations (25),(26),(27),(28)
 9: $\{\beta_y, \beta_{f_m}, \{\alpha_{mk}, \lambda_{mk}\}_k\}_m \leftarrow$
 10: M-STEP($\{\boldsymbol{\mu}_{w_m}, \Sigma_{w_m}, \boldsymbol{\mu}_{f_m}, \Sigma_{f_m}\}_m$)
 // pruning
 11: **for all** m **do**
 12: **if** $\alpha_{mk} > 1e3 \quad \forall k = 1, \dots, K$ **then**
 13: $M \leftarrow M - 1; C \leftarrow C \setminus c_m$
 14: **end if**
 15: **end for**
 16: **end for**

4 Extension to Finitely Many Local Nonparametric Models

An interesting component of local Gaussian regression is that it easily extends to a model with finitely many local nonparametric, Gaussian process models. Marginalising out the weights \mathbf{w}_m in Eq. (17) gives the marginal

$$p(\mathbf{y}, \mathbf{f} | X, \theta) = \mathcal{N}(\mathbf{y}; \mathbf{1}^T \mathbf{f}, \beta_y^{-1} I_N) \prod_{m=1}^M \mathcal{N}(\mathbf{f}_m; 0, C_m) \quad (32)$$

where $p(f_m) = \mathcal{N}(\mathbf{f}_m; 0, C_m)$ is a Gaussian process prior over function f_m with a (degenerate) covariance function $C_m = \beta_{f_m}^{-1} I_N + \phi_m^T A_m^{-1} \phi_m$. Replacing the finitely many local features ξ_m with infinitely many features results in the local nonparametric model with covariance function $\kappa_m(x, x') = \beta_{f_m}^{-1} \delta_{ij} + \eta_m(x) \hat{\kappa}_m(x, x') \eta_m(x')$. The exact posterior over $\mathbf{f}_x := \mathbf{f}(x)$ is

$$p(\mathbf{f}_x | \mathbf{y}) = \mathcal{N}(\mathbf{f}_x; \boldsymbol{\mu}_x, \Sigma_{xx}) \quad (33)$$

with

$$\Sigma_{xx}^{mr} = \text{cov}(f^m(x), f^r(x')) \quad (34)$$

$$= \delta_{mr} k_{xx'}^m - k_{xX}^m \left(\sum_s k_{XX}^s + \beta_y^{-1} \mathbf{I} \right)^{-1} k_{Xx'}^r \quad (35)$$

$$\boldsymbol{\mu}_x^m = k_{xX}^m \left(\sum_s k_{XX}^s + \beta_y^{-1} \mathbf{I} \right)^{-1} \mathbf{y} \quad (36)$$

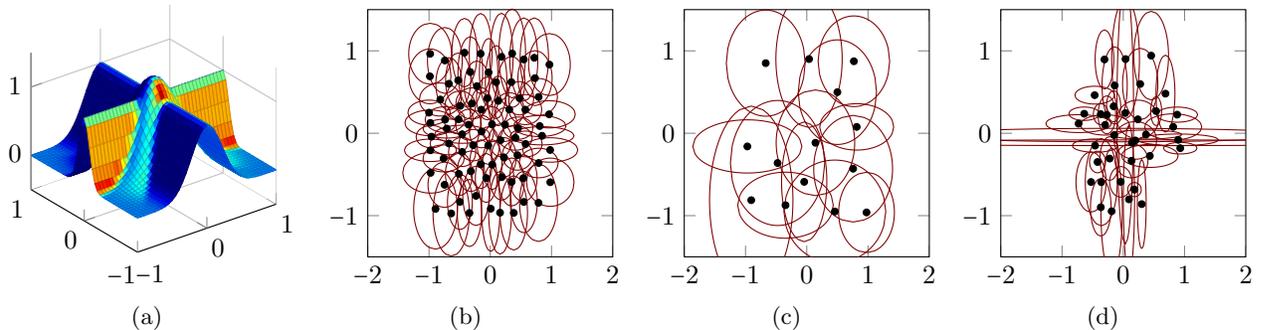


Figure 3: (a) 2D cross function, local models learnt by (b) LWPR, (c) LGR and (d) ME .

Table 1: Predictive performance on Cross Function

	nMSE w/o LSL	opt w_{gen}	# of LMs	nMSE with LSL	opt w_{gen}	# of LMs
LWPR	0.234	0.2	13	0.0365	0.1	45.5
GLR	0.069	1.0	23.8	0.0137	0.9	23.3
ME	0.169	0.5	21.2	0.0313	0.4	62.2

So computing the posterior covariance Σ_{xx}^{mr} requires one inversion of an $N \times N$ matrix. It remains to be seen to what extent variational bounds on the parametric forms presented in the previous sections can be transported to a *local nonparametric* Gaussian regression framework.

5 Experiments

We evaluate and compare to mixture of experts and locally weighted projection regression (LWPR) – an extension of LWR suitable for regression in high dimensional space [17] – on two data sets. The first is data drawn from the “cross function” in Figure 3, often used to demonstrate locally weighted learning. For the second comparison we learn inverse dynamics of a SARCOS anthropomorphic robot arm with seven degrees of freedom [8].

In order to compare to mixture of experts we assume linear expert models. A prior of the form $\mathcal{N}(\mathbf{w}_m; 0, A_m^{-1})$ on each experts regression parameters and normalized gaussian kernels for the mixture components are used in our implementation to make it as comparable as possible. To compute the posterior in the E-step a mean field approximation is employed.

In both experiments, LWPR performed multiple cycles through the data sets. Both the local Gaussian regression and the mixture of experts implementation are executed based on Algorithm 1 and are allowed an additional 1000 iterations to reach convergence.

5.1 Data from the ‘Cross Function’

We used 2,000 uniformly distributed training inputs, with zero mean Gaussian noise of variance $(0.2)^2$. The test set is a regular grid of 1641 edges without noise and is used to evaluate how well the underlying function is captured. The initial length scale was set to $\lambda_m = 0.3, \forall m$, and we ran each method with and without lengthscale learning (LSL) for $w_{\text{gen}} = 0.1, 0.2, \dots, 1.0$. All results presented here are results averaged over 5 randomly seeded runs. Table 1 presents the top performance achieved by each method, including what the optimal setting for w_{gen} was and how many local models were used. In both settings (with or without LSL) LGR outperforms both LWPR and ME in terms of accuracy as well as number of local models used. To understand the role of parameter w_{gen} better, we also summarize the normalized mean squared error as a function of parameter w_{gen} for all 3 methods in Figure 4 with (right) and without (left) LSL. The key message of this graph is that the parameter w_{gen} does not affect the performance of LGR greatly. While accuracy slightly improves with increasing w_{gen} it is not dramatic. Thus for LGR w_{gen} can be thought of as a trade-off parameter, for smaller w_{gen} the algorithm has to consider less potential local models, for a slight performance decrease. For LWPR and ME this relationship to w_{gen} is not that clear. Furthermore, although LGR has to consider more local models for larger w_{gen} (up to 2000 for $w_{\text{gen}} = 1$), we only see a slight increase in the number of local models in the final model, indicating that the pruning mechanism works very well.

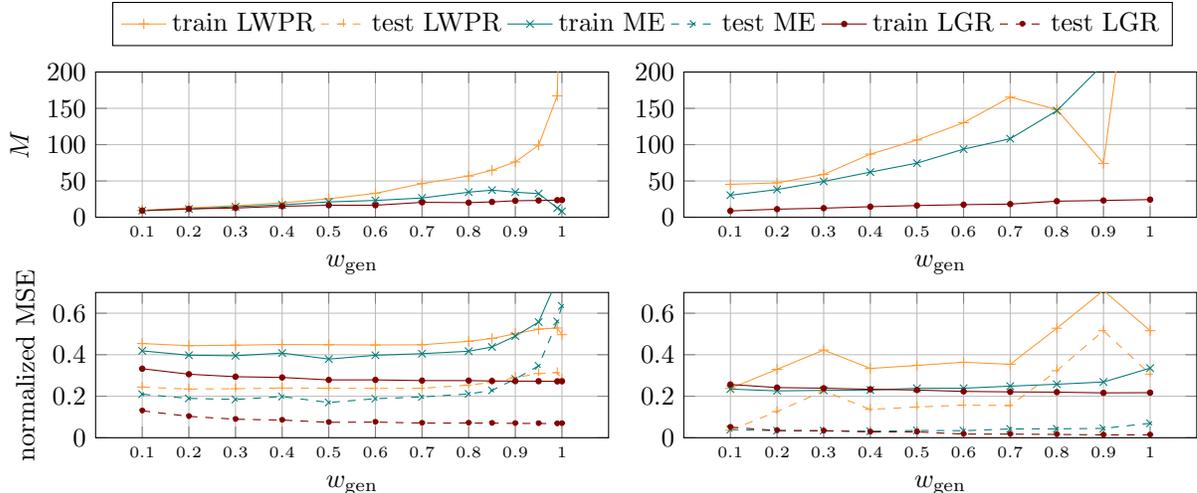


Figure 4: normalized mean squared error (nMSE) on 2D cross data without (left) and with (right) lengthscale learning.

Table 2: Predictive Performance on SARCOS task

	LWPR			LGR			ME		
	nMSE	(MSE)	# of LM	nMSE	(MSE)	# of LM	nMSE	(MSE)	# of LM
J1	0.045	(19.160)	419	0.015	(6.282)	366	0.071	(30.125)	249
J2	0.039	(8.938)	493	0.0195	(4.391)	361	0.0871	(19.564)	257
J3	0.034	(3.420)	483	0.014	(1.382)	359	0.074	(7.319)	243
J4	0.024	4.552	384	0.016	(3.118)	348	0.039	(7.348)	239
J5	0.064	(0.060)	514	0.035	(0.0336)	354	0.058	(0.0552)	438
J6	0.075	(0.221)	519	0.027	(0.0786)	359	0.124	(0.363)	234
J7	0.03	(0.203)	405	0.023	(0.1578)	358	0.0235	(0.1589)	383

Finally, we show representative results of the shape of learned local models for LWPR, LGR and ME In Figure 3, nicely illustrating the key difference between the three methods: In LWPR local models don't know of each other and thus aim to find the best linear approximation to the function. In both ME and LGR, the local models know of each other and collaborate to fit the function.

5.2 Inverse Dynamics Learning Task

The SARCOS data contains 44,484 training data points and 4,449 test data points. The 21 input variables represent joint positions, velocities and accelerations for the 7 joints. The task is to predict the 7 joint torques. In Table 2 we show the predictive performance of LWPR, ME and LGR when trained with lengthscale learning and with $w_{\text{gen}} = 0.3$. LGR outperforms LWPR and ME in terms of accuracy for almost all joints. However, the true advantage of LGR lies in the fact the number of hand tuned parameters is reduced to setting the learning rate for the gradient descent updates and setting the parameter w_{gen} .

6 Conclusion

We have taken a top-down approach to developing a probabilistic localised regression algorithm: We start with the generative model of Gaussian generalized linear regression, which amounts to a fully connected graph and thus has cubic inference cost. To break down the computational cost of inference, we first introduce the idea of localised feature functions as local models, which can be extended to local nonparametric models. In a second step, we argue that because of the localisation these local models are approximately independent. We exploit that fact through a variational approximation that reduced computational complexity to local computations. Empirical evaluation suggests that LGR successfully addresses the problem of fitting hyperparameters inherent in locally weighted regression. A final step left for future work is to re-formulate our algorithm into an incremental version that can deal with a continuous stream of incoming data.

References

- [1] Stefan Schaal and Christopher G Atkeson. Constructive incremental learning from only local information. *Neural Computation*, 10(8):2047–2084, 1998.
- [2] Léon Bottou and Vladimir Vapnik. Local learning algorithms. *Neural computation*, 4(6):888–900, 1992.
- [3] Robert A Jacobs, Michael I Jordan, Steven J Nowlan, and Geoffrey E Hinton. Adaptive mixtures of local experts. *Neural computation*, 3(1):79–87, 1991.
- [4] Steve Waterhouse, David MacKay, and Tony Robinson. Bayesian methods for mixtures of experts. *Advances in neural information processing systems*, pages 351–357, 1996.
- [5] Lauren Hannah, David M Blei, and Warren B Powell. Dirichlet process mixtures of generalized linear models. *Journal of Machine Learning Research*, 12:1923–1953, 2011.
- [6] Jo-Anne Ting, Mrinal Kalakrishnan, Sethu Vijayakumar, and Stefan Schaal. Bayesian Kernel Shaping for Learning Control. In *Neural information processing systems*, 2008.
- [7] Narayanan U Edakunni, Stefan Schaal, and Sethu Vijayakumar. Kernel carpentry for online regression using randomly varying coefficient model. In *Proceedings of the international joint conference on artificial intelligence (IJCAI)*, 2007.
- [8] C.E. Rasmussen and C.K.I. Williams. *Gaussian Processes for Machine Learning*. MIT Press, 2006.
- [9] C.A. Micchelli, Y. Xu, and H. Zhang. Universal kernels. *Journal of Machine Learning Research*, 7:2651–2667, 2006.
- [10] M. N Gibbs. *Bayesian Gaussian processes for regression and classification*. PhD thesis, University of Cambridge, 1997.
- [11] C.E. Rasmussen and Z. Ghahramani. Infinite mixtures of Gaussian process experts. In *Advances in neural information processing systems*. MIT Press, 2002.
- [12] R.M. Neal. *Bayesian learning for neural networks*. Springer Verlag, 1996.
- [13] M.E. Tipping. Sparse Bayesian learning and the relevance vector machine. *The Journal of Machine Learning Research*, 1:211–244, 2001.
- [14] Aaron D’Souza, Sethu Vijayakumar, and Stefan Schaal. The bayesian backfitting relevance vector machine. In *Proceedings of the twenty-first international conference on Machine learning*, page 31. ACM, 2004.
- [15] M.K. Titsias. Variational learning of inducing variables in sparse Gaussian processes. In *Artificial Intelligence & Statistics (AISTATS)*, volume 12, 2009.
- [16] Joaquin Quinonero-Candela and Ole Winther. Incremental gaussian processes. In *Advances in neural information processing systems*, pages 1001–1008, 2002.
- [17] Sethu Vijayakumar and Stefan Schaal. Locally weighted projection regression: An $o(n)$ algorithm for incremental real time learning in high dimensional space. In *Proceedings of the Seventeenth International Conference on Machine Learning (ICML 2000)*, volume 1, pages 288–293, 2000.