

Generalization Errors and Learning Curves for Regression with Multi-task Gaussian Processes

Authored by:

Kian M. Chai

Abstract

We provide some insights into how task correlations in multi-task Gaussian process (GP) regression affect the generalization error and the learning curve. We analyze the asymmetric two-task case, where a secondary task is to help the learning of a primary task. Within this setting, we give bounds on the generalization error and the learning curve of the primary task. Our approach admits intuitive understandings of the multi-task GP by relating it to single-task GPs. For the case of one-dimensional input-space under optimal sampling with data only for the secondary task, the limitations of multi-task GP can be quantified explicitly.

1 Paper Body

Gaussian processes (GPs) (see e.g., [1]) have been applied to many practical problems. In recent years, a number of models for multi-task learning with GPs have been proposed to allow different tasks to leverage on one another [2?5]. While it is generally assumed that learning multiple tasks together is beneficial, we are not aware of any work that quantifies such benefits, other than PACbased theoretical analysis for multi-task learning [6?8]. Following the tradition of the theoretical works on GPs in machine learning, our goal is to quantify the benefits using average-case analysis. We concentrate on the asymmetric two-tasks case, where the secondary task is to help the learning of the primary task. Within this setting, the main parameters are (1) the degree of ?relatedness? ? between the two tasks, and (2) the ratio ?S of total training data for the secondary task. While higher —?— and lower ?S is clearly more beneficial to the primary task, the extent and manner that this is so has not been clear. To address this, we measure the benefits using generalization error, learning curve and optimal error, and investigate the influence of ? and ?S on these quantities. We will give non-trivial lower and upper bounds on the generalization error and the learning curve. Both types of bounds are important in providing assurance on the quality of predictions: an upper bound provides an estimate of the amount of training data needed to attain a minimum performance level, while a lower bound provides an understanding of the limitations

of the model [9]. Our approach relates multi-task GPs to single-task GPs and admits intuitive understandings of multi-task GPs. For one-dimensional input-space under optimal sampling with data only for the secondary task, we show the limit to which error for the primary task can be reduced. This dispels any misconception that abundant data for the secondary task can remedy no data for the primary task.

2.1

Preliminaries and problem statement Multi-task GP regression model and setup

The multi-task Gaussian process regression model in [5] learns M related functions $\{f_m\}_{m=1}^M$ by placing a zero mean GP prior which directly induces correlations between tasks. Let y_m be an 1

observation of the m th function at x . Then the model is given by
$$y_m \sim \mathcal{N}(f_m(x), \sigma_m^2), \quad (1)$$

where k_x is a covariance function over inputs, and K_f is a positive semi-definite matrix of inter-task similarities, and σ_m^2 is the noise variance for the m th task. The current focus is on the two tasks case, where the secondary task S is to help improve the performance of the primary task T ; this is the asymmetric multi-task learning as coined in [10]. We fix K_f to be a correlation matrix, and let the variance be explained fully by k_x (the converse has been done in [5]). Thus K_f is fully specified by the correlation $\rho \in [0, 1]$ between the two tasks. We further fix the noise variances of the two tasks to be the same, say σ^2 . For the training data, there are n_T (resp. n_S) observations at locations X_T (resp. X_S) for task T (resp. S). We use $n = n_T + n_S$ for the total number of observations, $\rho_S = n_S/n$ for the proportion of observations for task S , and also $X = X_T \cup X_S$. The aim is to infer the noise-free response f_T for task T at x^* . See Figure 1. The covariance matrix of the noisy training data is $K(\rho) + \sigma^2 I$, where
$$K(\rho) = \begin{bmatrix} K_{TT} & K_{TS} \\ K_{ST} & K_{SS} \end{bmatrix} \quad (2)$$

K_{TT} and K_{TS} (resp. K_{SS}) is the matrix of covariances (due to k_x) between locations in X_T (resp. X_S); K_{TS} is the matrix of cross-covariances from locations in X_T to locations in X_S ; and K_{ST} is K_{TS} transposed. The posterior variance at x^* for task T is

$$v_T(x^*) = \sigma^2 \left(\frac{1}{n} + \frac{k(x^*, X_T)}{K_{TT} + \sigma^2 I} \right), \quad (3)$$

and $k(x^*, X_T)$ (resp. $k(x^*, X_S)$) is the vector of covariances (due to k_x) between locations in X_T (resp. X_S) and x^* . Where appropriate and clear from context, we will suppress some of the parameters in $v_T(x^*)$, or use X for (X_T, X_S) . Note that $v_T(x^*) = v_T(x^*)$, so that $v_T(x^*)$ is the same as $v_T(x^*)$; for brevity, we only write the former. If the GP prior is correctly specified, then the posterior variance (3) is also the generalization error at x^* [1, 7.3]. The latter is defined as $h(f_T(x^*)) = \mathbb{E} \|f_T(x^*) - \hat{f}_T(x^*)\|^2$, where $\hat{f}_T(x^*)$ is the posterior mean at x^* for

task T , and the expectation is taken over the distribution from which the true function f_T is drawn. In this paper, in order to distinguish succinctly from the generalization error introduced in the next section, we use posterior variance to mean the generalization error at x^* . Note that the actual y -values observed at X do not effect the posterior variance at any test location. Problem statement Given the above setting, the aim is to investigate how training observations for task S can benefit the predictions for task T . We measure the benefits using generalization error, learning curve and optimal error, and investigate how these quantities vary with ρ and γ_S . 2.2

Generalization errors, learning curves and optimal errors

We outline the general approach to obtain the generalization error and the learning curve [1, §7.3] under our setting, where we have two tasks and are concerned with the primary task T . Let $p(x)$ be the probability density, common to both tasks, from which test and training locations are drawn, and assume that the GP prior is correctly specified. The generalization error for task T is obtained by averaging the posterior variance for task T over x^* , and the learning curve for task T is obtained by averaging the generalization error over training sets X :
$$\text{generalization error: } T(\gamma, \gamma_{n2}, X_T, X_S) \stackrel{\text{def}}{=} \int T^2(x^*, \gamma, \gamma_{n2}, X_T, X_S) p(x^*) dx^* \quad (4)$$

$$\text{learning curve: } T(\gamma, \gamma_n, \gamma_S, n) = \int T(\gamma, \gamma_n, X_T, X_S) p(X) dX, \quad (5)$$
 where the training locations in X are drawn i.i.d, that is, $p(X)$ factorizes completely into a product of $p(x)$ s. Besides averaging T to obtain the learning curve, one may also use the optimal experimental design methodology and minimize T over X to find the optimal generalization error [11, chap. II]:
$$\text{optimal error: } T(0, \gamma_{n2}, X_T, X_S) \stackrel{\text{def}}{=} \min_X T(\gamma, \gamma_n, X_T, X_S).$$

$$\text{opt } T(\gamma, \gamma_n, \gamma_S, n) = \min_X T(\gamma, \gamma_n, X_T, X_S) \quad (6)$$

Both can reduce to single-task GP cases; the former discards training observations at X_S , while the latter includes them. Similar analogues to single-task GP avg opt opt cases for $\text{avg } T(0, \gamma_n, \gamma_S, n)$ and $T(1, \gamma_n, \gamma_S, n)$, and $T(0, \gamma_n, \gamma_S, n)$ and $T(1, \gamma_n, \gamma_S, n)$ can be avg opt obtained. Note that T and T are well-defined since $\gamma_S n = n_S \neq 0$ by the definition of γ_S . 2

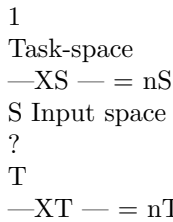
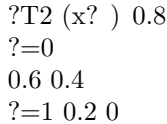


Figure 1: The two tasks S and T have task correlation ρ . The data set X_T (resp. X_S) for task T (resp. S) consists of the γ s (resp. s). The test location x^* for task T is denoted by γ^* . 2.3



0
0.2
0.4
x? 0.6
0.8
1

Figure 2: The posterior variances of each test location within $[0, 1]$ given data \mathbf{y}_T at $1/3$ and $2/3$ for task T , and \mathbf{y}_S at $1/5$, $1/2$ and $4/5$ for task S .

Eigen-analysis

We now state known results of eigen-analysis used in this paper. Let λ_i and $\phi_i(x)$ be the eigenvalues and eigenfunctions of the covariance function $k(x, x')$ under the measure $p(x)dx$: they satisfy the integral equation $\int_0^1 k(x, x') \phi_i(x') p(x') dx' = \lambda_i \phi_i(x)$. Let $\mathbf{x} = (x_1, \dots, x_n)$ be the locations in \mathcal{X} sampled, $\lambda_i = \lambda_i(\mathbf{x})$ be the eigenvalues of $\mathbf{K}(\mathbf{x}) = \int_0^1 k(\mathbf{x}, \mathbf{x}') p(\mathbf{x}') d\mathbf{x}'$ from $p(\mathbf{x})$, then $\lambda_i = \lim_{n \rightarrow \infty} \lambda_i(\mathbf{x})$. However, for finite n used in practice, the estimate λ_i/n for λ_i is better for the larger eigenvalues than for the smaller ones. Additionally, in one-dimension with uniform $p(x)$ on the unit interval, if $k(x, x')$ satisfies the Sacks-Ylvisaker conditions of order r , then $\lambda_i \sim i^{-2r-2}$ in the limit $i \rightarrow \infty$ [11, Proposition IV.10, Remark IV.2]. Broadly speaking, an order r process is exactly r times mean square differentiable. For example, the stationary Ornstein-Uhlenbeck process is of order $r = 0$.

3

Generalization error

In this section, we derive expressions for the generalization error (and the bounds thereon) for the two-tasks case in terms of the single-task one. To illustrate and further motivate the problem, Figure 2 plots the posterior variance $\sigma^2_T(x, \mathbf{y}_T, \mathbf{y}_S)$ as a function of x given two observations for task T and three observations for task S . We roughly follow [13, Fig. 2], and use squared exponential covariance function with length-scale 0.11 and noise variance $\sigma^2 = 0.05$. Six solid curves are plotted, corresponding, from top to bottom, to $\sigma^2 = 0, 1/8, 1/4, 1/2, 3/4$ and 1. The two dashed curves enveloping each solid curve are the lower and upper bounds derived in this section; the dashed curves are hardly visible because the bounds are rather tight. The dotted line is the prior noise variance. Similar to the case of single-task learning, each training point creates a depression on the $\sigma^2_T(x, \mathbf{y}_T, \mathbf{y}_S)$ surface [9, 13]. However, while each training point for task T creates a ‘full’ depression that reaches the prior noise variance (horizontal dotted line at 0.05), the depression created by each training point for task S depends on σ^2 , ‘deeper’ depressions for larger σ^2 . From the figure, and also from definition, it is clear that the following trivial bounds on $\sigma^2_T(x, \mathbf{y}_T, \mathbf{y}_S)$ hold: Proposition 1. For all $x \in [0, 1]$, $\sigma^2_T(x, \mathbf{y}_T, 1) \leq \sigma^2_T(x, \mathbf{y}_T, \mathbf{y}_S) \leq \sigma^2_T(x, \mathbf{y}_T, 0)$. Integrating wrt to x then gives the following corollary: Corollary 2. $\int_0^1 \sigma^2_T(x, \mathbf{y}_T, \mathbf{y}_S) dx \leq \int_0^1 \sigma^2_T(x, \mathbf{y}_T, 1) dx \leq \int_0^1 \sigma^2_T(x, \mathbf{y}_T, 0) dx$. Sections 3.2 and 3.3 derive lower and upper bounds that are tighter than the above trivial bounds. Prior to the bounds, we consider a degenerate case to illustrate the limitations of multi-task learning. 3.1

The degenerate case of no training data for primary task

It is clear that if there is no training data for the secondary task, that is, if $X_S = \emptyset$, then $\mathbb{E}T_2(x; 1) = \mathbb{E}T_2(x; \cdot, \cdot) = \mathbb{E}T_2(x; 0)$ for all x and \cdot . In the converse case where there is no training data for the primary task, that is, $X_T = \emptyset$, we instead have the following proposition: 3

Proposition 3. For all x , $\mathbb{E}T_2(x; \cdot, \cdot, \cdot, X_S) = \mathbb{E}T_2(x; 1, \cdot, X_S) + (1 - \mathbb{E}T_2)k$. Proof.

$$\mathbb{E}_x T_2(x; \cdot, \cdot, \cdot, X_S) = k \mathbb{E}T_2(kxS)T(KSS + \mathbb{E}I) + kxS$$

$$x = (1 - \mathbb{E}T_2)k + \mathbb{E}T_2 k \mathbb{E}T_2(kxS)T(KSS + \mathbb{E}I) + kxS$$

$= (1 - \mathbb{E}T_2)k + \mathbb{E}T_2 \mathbb{E}T_2(x; 1, \cdot, X_S)$. Hence the posterior variance is a weighted average of the prior variance k and the posterior variance at perfect correlation. When the cardinality of X_S increases under infill asymptotics [14, §3.3], $\lim_{n \rightarrow \infty} \mathbb{E}T_2(x; 1, \cdot, X_S) = 0$

$=$

$$\lim_{n \rightarrow \infty} \mathbb{E}T_2(x; \cdot, \cdot, \cdot, X_S) = (1 - \mathbb{E}T_2)k.$$

(7)

This is the limit for the posterior variance at any test location for task T , if one has training data only for the secondary task S . This is because a correlation of $\mathbb{E}T_2$ between the tasks prevents any training location for task S from having correlation higher than $\mathbb{E}T_2$ with a test location for task T . Suppose correlations in the input-space are given by an isotropic covariance function $k(x) = k(\|x - x_0\|)$. If we translate correlations into distances between data locations, then any training location from task S is beyond a certain radius from any test location for task T . In contrast, a training location from task T may lay arbitrarily close to a test location for task T , subject to the constraints of noise. We obtain the generalization error in this degenerate case, by integrating Proposition 3 wrt $p(x)dx$ and using the fact that the mean prior variance is given by the sum of the process eigenvalues. P? Corollary 4. $T(\cdot, \mathbb{E}I, \cdot, X_S) = \mathbb{E}T_2(1, \mathbb{E}I, \cdot, X_S) + (1 - \mathbb{E}T_2) \sum_{i=1}^{\infty} \lambda_i$. 3.2

A lower bound

When $X_T \neq \emptyset$, the correlations between locations in X_T and locations in X_S complicate the situation. However, since $\mathbb{E}T_2(\cdot)$ is a continuous and monotonically decreasing function of \cdot , there exists an $\alpha \in [0, 1]$, which depends on \cdot, x and X , such that $\mathbb{E}T_2(\cdot) = \alpha \mathbb{E}T_2(1) + (1 - \alpha) \mathbb{E}T_2(0)$. That α depends on x obstructs further analysis. The next proposition gives a lower bound $\mathbb{E}T_2(\cdot)$ of the $\mathbb{E}T_2$ of x in same form satisfying $\mathbb{E}T_2(1) \leq \mathbb{E}T_2(\cdot) \leq \mathbb{E}T_2(0)$, where the mixing proportion is independent of x def Proposition 5. Let $\mathbb{E}T(x; \cdot, \cdot) = \alpha \mathbb{E}T(x; 1, \cdot) + (1 - \alpha) \mathbb{E}T(x; 0, \cdot)$. Then for all x : (a) $\mathbb{E}T_2(x; \cdot, \cdot) \geq \mathbb{E}T_2(x; \cdot, \cdot)$ (b) $\mathbb{E}T_2(x; \cdot, \cdot) \geq \mathbb{E}T_2(x; \cdot, \cdot) \geq \mathbb{E}T_2(x; 0, \cdot) \geq \mathbb{E}T_2(x; \cdot, 1)$

(c) $\arg \max_{\alpha} \mathbb{E}T_2(x; \cdot, \cdot) \geq \mathbb{E}T_2(x; \cdot, \cdot) \geq 1/2$. The proofs are in supplementary material §S.2. The lower bound $\mathbb{E}T_2(\cdot)$ depends explicitly on $\mathbb{E}T_2$ for task S , through the gap. It depends implicitly on $\mathbb{E}S$, which is the proportion of observations between $\mathbb{E}T_2(1)$ and $\mathbb{E}T_2(0)$. If there is no training data for the primary task, i.e., if $\mathbb{E}S = 1$, the bound reduces to Proposition 3, and becomes exact for all values of \cdot . If $\mathbb{E}S = 0$, the bound is also exact. For

$\beta \in \{0, 1\}$, the bound is exact when $\beta \in \{1/2, 0, 1\}$. As from Figure 2 and later from our simulation results in section 5.3, this bound is rather tight. Part (b) of the proposition states the tightness of the bound: it is no more than factor β of the gap between the trivial bounds $\sigma_T^2(0)$ and $\sigma_T^2(1)$. Part (c) of the proposition says that the bound is least tight for a value of β greater than $1/2$. We provide an intuition on Proposition 5a. Let f_T^1 (resp. f_T^0) be the posterior mean of the single-task GP when $\beta = 1$ (resp. $\beta = 0$). Contrasted with the multi-task predictor f_T , f_T^1 directly involves the noisy observations for task T at X_S , so it has more information on task T . Hence, predicting $f_T^1(x^*)$ gives the trivial lower bound $\sigma_T^2(1)$ on $\sigma_T^2(\beta)$. The tighter bound $\sigma_T^2(\beta)$ is obtained by ‘throwing β away’ information and predicting $f_T^1(x^*)$ with probability β and $f_T^0(x^*)$ with probability $(1 - \beta)$. Finally, the next corollary is readily obtained from Proposition 5a by integrating wrt $p(x^*)dx^*$. This is possible because β is independent of x^* . Corollary 6. Let $T(\beta, \beta n_2, X_T, X_S)$ def $= \beta T(1, \beta n_2, X_T, X_S) + (1 - \beta) T(0, \beta n_2, X_T, X_S)$. Then $\sigma_T^2(\beta, \beta n_2, X_T, X_S) \leq T(\beta, \beta n_2, X_T, X_S)$. \square 3.3

An upper bound via equivalent isotropic noise at X_S

The following question motivates our upper bound: if the training locations in X_S had been observed for task T rather than for task S , what is the variance σ_T^2 of the equivalent isotropic noise at X_S so that

that the posterior variance remains the same? To answer this question, we first refine the definition of $\sigma_T^2(\beta)$ to include a different noise variance parameter s_2 for the X_S observations: $\sigma_T^2(\beta, \beta n_2, s_2, X_T, X_S)$ def $= k(\beta, \beta n_2, s_2, X_T, X_S)$ (8) $= k(\beta, \beta n_2, s_2, X_T, X_S) + 0n_2 s_2$ I cf. (3). We may suppress the parameters β, X_T and X_S when writing $\sigma_T^2(\beta)$. The variance $\sigma_T^2(\beta, \beta n_2, s_2, X_T, X_S)$ of the equivalent isotropic noise is a function of β defined by the equation $\sigma_T^2(\beta, \beta n_2, s_2, X_T, X_S) = \sigma_T^2(\beta, \beta n_2, s_2, X_T, X_S)$. (9) For any β there is always a β that satisfies the equation because the difference $\sigma_T^2(\beta, \beta n_2, s_2, X_T, X_S) - \sigma_T^2(\beta, \beta n_2, s_2, X_T, X_S)$ is a continuous and monotonically decreasing function of β . To make progress, we seek an upper bound $\sigma_T^2(\beta, \beta n_2, s_2, X_T, X_S)$ that is independent of the choice of β : $\sigma_T^2(\beta, \beta n_2, s_2, X_T, X_S) \leq \sigma_T^2(\beta, \beta n_2, s_2, X_T, X_S)$ for all test locations. Of $\sigma_T^2(\beta, \beta n_2, s_2, X_T, X_S)$, which is the minimum possible σ_T^2 interest is the tight upper bound $\sigma_T^2(\beta, \beta n_2, s_2, X_T, X_S)$, given in the next proposition. $\beta + \beta n_2 + \beta n_2$ be the maximum eigenvalue of $K(x, \beta)$ def $\sigma_T^2(\beta, \beta n_2, s_2, X_T, X_S)$ Proposition 7. Let $\beta = \beta(\beta = \beta n_2 + \beta n_2)$ and $\beta n_2 = \beta n_2$ SS βn_2 . The bound is tight in this sense: for any β Then for all β , $\sigma_T^2(\beta, \beta n_2, s_2, X_T, X_S) \leq \sigma_T^2(\beta, \beta n_2, s_2, X_T, X_S) \leq \sigma_T^2(\beta, \beta n_2, s_2, X_T, X_S)$ if β $\sigma_T^2(\beta, \beta n_2, s_2, X_T, X_S) \leq \sigma_T^2(\beta, \beta n_2, s_2, X_T, X_S)$, then β $\sigma_T^2(\beta, \beta n_2, s_2, X_T, X_S) \leq \sigma_T^2(\beta, \beta n_2, s_2, X_T, X_S)$. Proof sketch. Matrix $K(\beta)$ may be factorized as

$$K(\beta) = K_T^T I_0 K(\beta) = x_0^T I_0 K^T$$

$$K^T x_S \beta x^T K^T S$$

$$I_0$$

$$0. \beta I_0 I_0 \beta I$$

By using this factorization in the posterior variance (8) and taking out the $(kx^*)^T$

$$\text{def}$$

$$\sigma_T^2(\beta, \beta n_2, s_2)$$

does not depend on γ_S even for moderate sizes n_S . Therefore, the lower bound is not as useful as the upper bound. Finally, if we refine T as we have done for T_2 in (8), we obtain the following corollary: γ_{2n, X_T, X_S} . Then Corollary 8. Let $T(\gamma, \gamma_{2n}, \gamma_{2n}, X_T, X_S) \stackrel{\text{def}}{=} T(1, \gamma_{2n}, \gamma T(\gamma, \gamma_{2n}, \gamma_{2n}, X_T, X_S)) \leq T(\gamma, \gamma_{2n}, \gamma_{2n}, X_T, X_S)$. 5

3.4

Exact computation of generalization error

The factorization of T_2 expressed by (12) allows the generalization error to be computed exactly in certain cases. We replace the quadratic form in (12) by matrix trace and then integrate out x to give P

$T(\gamma, \gamma_{2n}, X_T, X_S) = \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \text{tr} \left(\frac{1}{n} \sum_{j=1}^n \mathbb{E} [k(x_i, x_j)] M \right) \right]$, where \mathbb{E} denotes $\mathbb{E}(\gamma, \gamma_{2n}, \gamma_{2n})$, the expectations are taken over x , and M is an n -by- n matrix with $R_{P^2} M_{pq} \stackrel{\text{def}}{=} k(x_p, x_q) - k(x_p, x_q) p(x_q) \int p(x) dx = \sum_{i=1}^n \gamma_{2i} \gamma_i (x_p) \gamma_i (x_q)$, where $x_p, x_q \in X$. When the eigenfunctions $\gamma_i(\cdot)$ s are not bounded, the infinite-summation expression for M_{pq} is often difficult to use. Nevertheless, analytical results for M_{pq} are still possible in some cases using the integral expression. An example is the case of the squared exponential covariance function with normally distributed x , when the integrand is a product of three Gaussians.

4

Optimal error for the degenerate case of no training data for primary task

If training examples are provided only for task S , then task T has the following optimal performance. Proposition 9. Under optimal sampling on a 1-d space, if the covariance function satisfies Sacks $\gamma_{(2r+1)/(2r+2)}^2$ Ylvisaker conditions of order r , then $\text{opt}(\gamma) + (1 - \gamma_{2n}) \sum_{i=1}^n \gamma_i = T(\gamma, \gamma, 1, n) = \gamma(n_S P^2)^{2/2} \text{opt}^2$ Proof. We obtain $\text{opt} T(\gamma, \gamma, 1, n) = \gamma T(1, \gamma_n, 1, n) + (1 - \gamma_{2n}) \sum_{i=1}^n \gamma_i$ by minimizing Corollary 4 wrt X_S . Under the same conditions as the proposition, the optimal generalization error using the single-task GP decays with training set size n as $\gamma(n^{(2r+1)/(2r+2)})$ [11, Proposition V.3]. Thus $\gamma_{(2r+1)/(2r+2)}^2 \gamma_{(2r+1)/(2r+2)}^2 = \gamma(n_S)$. $\gamma^2 \text{opt} T(1, \gamma_n, 1, n) = \gamma \gamma(n_S P^2)$ A directly corollary of the above result is that one cannot expect to do better than $(1 - \gamma_{2n}) \sum_{i=1}^n \gamma_i$ on the average. As this is a lower bound, the same can be said for incorrectly specified GP priors.

5

Theoretical bounds on learning curve

Using the results from section 3, lower and upper bounds on the learning curve may be computed by averaging over the choice of X using Monte Carlo approximation.¹ For example, using Corollary 2 and integrating wrt $p(X)dX$ gives the following trivial bounds on the learning curve: $\text{avg} \text{avg}^2 \text{Corollary 10. } \text{avg} T(1, \gamma_n, \gamma_S, n) \leq T(\gamma, \gamma_n, \gamma_S, n) \leq T(0, \gamma_n, \gamma_S, n)$. The gap between the trivial bounds can be analyzed as follows. Recall that $\gamma_S n \leq N_0$ by definition, $\text{avg} \text{avg}^2 \text{so that } \text{avg} T(1, \gamma_n, \gamma_S, (1 - \gamma_S)n) = T(0, \gamma_n, \gamma_S, n)$. Therefore $T(1, \gamma_n, \gamma_S, n)$ is equivalent to $\text{avg}^2 T(0, \gamma_n, \gamma_S, n)$ scaled along the n -axis by the factor $(1 - \gamma_S) \in [0, 1]$, and hence the gap between the trivial bounds becomes wider with γ_S . In the rest of this section, we derive non-trivial theoretical bounds on the learning curve before

providing simulation results. Theoretical bounds are particularly attractive for high-dimensional input-spaces, on which Monte Carlo approximation is harder.

5.1

Lower bound

For the single-task GP, a lower bound on its learning curve is $\sum_{i=1}^{n_2} \frac{\sigma_i^2}{(\sigma^2 + \sigma_i^2)}$ [15]. We shall call this the single-task OV bound. This lower bound can be combined with Corollary 6. Proposition 11. $\text{avg}(\sigma^2, \sigma^2, \sigma^2, n) \leq \sigma^2 + (1 - \sigma^2) \frac{\sigma^2}{\sigma^2 + \sigma^2} \sum_{i=1}^n \frac{\sigma_i^2}{\sigma^2 + \sigma_i^2} + (1 - \sigma^2) \frac{\sigma^2}{\sigma^2 + \sigma^2} \sum_{i=1}^n \frac{\sigma_i^2}{\sigma^2 + \sigma_i^2}$ or equivalently, $\text{avg} T(\sigma^2, \sigma_n^2, \sigma_S^2, n) \leq \sigma^2$

$\sum_{i=1}^n \frac{\sigma_i^2}{\sigma^2 + \sigma_i^2}$

or equivalently, $\text{avg} T(\sigma^2, \sigma_n^2, \sigma_S^2, n) \leq \sigma^2$

$\sum_{i=1}^n \frac{\sigma_i^2}{\sigma^2 + \sigma_i^2}$

1

$b_i \sigma_i^2, \sigma^2 + \sigma_i^2$

with $b_i \text{ def } =$

$\sigma^2 + (1 - \sigma^2) \frac{\sigma_S^2}{\sigma^2 + \sigma_S^2} \sigma_i^2, \sigma^2 + (1 - \sigma^2) \frac{\sigma_S^2}{\sigma^2 + \sigma_S^2} \sigma_i^2$

$\frac{\sigma^2}{\sigma^2 + \sigma_i^2} \sigma_i^2 \text{ def } \sigma^2 + (1 - \sigma^2) \frac{\sigma_S^2}{\sigma^2 + \sigma_S^2} \sigma_i^2$, with $b_i = \sigma^2 + (1 - \sigma^2) \frac{\sigma_S^2}{\sigma^2 + \sigma_S^2} \sigma_i^2 + \sigma_i^2$

Approximate lower bounds are also possible, by combining Corollary 6 and approximations in, e.g., [13].

6

Proof sketch. To obtain the first inequality, we integrate Corollary 6 wrt to $p(X)dX$, and apply the single-task OV bound twice. For the second inequality, its i th summand is obtained by combining the corresponding pair of i th summands in the first inequality. The third inequality is obtained from the second by swapping the denominator of b_i with that of $\sigma_i^2 / (\sigma^2 + \sigma_i^2)$ for every i . For fixed σ^2, σ_S^2 and n , denote the above bound by OV^* . Then OV^0 and OV^1 are both single task bounds. In particular, from Corollary 10, we have that the OV^1 is a lower bound on $\text{avg} T(\sigma^2, \sigma_n^2, \sigma_S^2, n)$. From the first expression of the above proposition, it is clear from the ‘mixture’ nature of the bound that the two-tasks bound OV^* is always better than OV^1 . As σ^2 decreases, the two-tasks bound moves towards the OV^0 ; and as σ_S^2 increases, the gap between OV^0 and OV^1 increases. In addition, the gap is also larger for rougher processes, which are harder to learn. Therefore, the relative tightness of OV^* over OV^1 is more noticeable for lower σ^2 , higher σ_S^2 and rougher processes. The second expression in the Proposition 11 is useful for comparing with the OV^1 . Each summand for the two-tasks case is a factor b_i of the corresponding summand for the single-task case. Since $b_i \in [1, (1 - \sigma^2) \frac{\sigma_S^2}{\sigma^2 + \sigma_S^2} / (1 - \sigma^2) \frac{\sigma_S^2}{\sigma^2 + \sigma_S^2}]$, OV^* is more than OV^1 by at most $(1 - \sigma^2) \frac{\sigma_S^2}{\sigma^2 + \sigma_S^2} / (1 - \sigma^2) \frac{\sigma_S^2}{\sigma^2 + \sigma_S^2}$ times. Similarly, the third expression of the proposition is useful for comparing with OV^0 : each summand for the the two-tasks case is a factor $b_{0i} \in [(1 - \sigma^2) \frac{\sigma_S^2}{\sigma^2 + \sigma_S^2}, 1]$ of the corresponding single-task one. Hence, OV^* is less than OV^0 by up to $(1 - \sigma^2) \frac{\sigma_S^2}{\sigma^2 + \sigma_S^2}$ times. In terms of the lower bound, this is the limit to which multi-task learning can outperform the single-task learning that ignores the secondary task. 5.2

Upper bound using equivalent noise

0.6
0.6
0.4
0.4
0.2
0.2
n 0
50
100
150
200
250
n
300
0
(a) $\gamma_2 = 1/2, \gamma_S = 1/2$
50
100
150
200
250
300
(b) $\gamma_2 = 3/4, \gamma_S = 3/4$

Figure 3: Comparison of various bounds for two settings of (γ, γ_S) . Each graph plots $\text{avg } T$ against n and consists of the ‘true’ multi-task learning curve (middle), the theoretical lower/upper bounds of Propositions 11/13 (lower/upper), the empirical trivial lower/upper bounds using Corollary 10 (lower/upper), and the empirical lower/upper bounds using Corollaries 6/8 ($\gamma/4$). The thickness of the ‘true’ multi-task learning curve reflects 95% confidence interval. unit variance squared exponential $k(x, x_0) = \exp[-(x - x_0)^2 / (2l^2)]$ with length-scale $l = 0.01$, the observation noise variance is $\sigma^2 = 0.05$, and the learning curves are computed for up to $n = 300$ training data points. When required, the average over x is computed analytically (see section 3.4). The empirical average over $X \stackrel{\text{def}}{=} XT - XS$, denoted by $\text{hh}T$, is computed over 100 randomly sampled training sets. The process eigenvalues λ_i needed to compute the theoretical bounds are given in [17]. Supplementary material S.6 gives further details. Learning curves for pairwise combinations of $\gamma \in \{1/8, 1/4, 1/2, 3/4\}$ and $\gamma_S \in \{1/4, 1/2, 3/4\}$ are computed. We compare the following: (a) the ‘true’ multi-task learning curve $\text{hh}T$ obtained by averaging T_2 over x and X ; (b) the theoretical bounds OV and FWO of Propositions 11 and 13; (c) the trivial upper and lower bounds that are single-task learning curves $\text{hh}T(0)$ and $\text{hh}T(1)$ obtained by averaging $T_2(0)$ and $T_2(1)$; and (d) the empirical lower bound $\text{hh}T$ and upper bound $\text{hh}T$ using Corollaries 6 and 8. Figure 3 gives some indicative plots of the curves. We summarize with the following observations: (a) The gap between the trivial bounds $\text{hh}T(0)$ and $\text{hh}T(1)$ increases with γ_S , as described at the start of section 5. (b) We

find the lower bound $\text{hhT}(\gamma)$ is a rather close approximation to the multi-task learning curve $\text{hhT}(\gamma)$, as evidenced by the much overlap between the lines and the middle lines in Figure 3. (c) The curve for the empirical upper bound $\text{hhT}(\gamma)$ using the equivalent noise method has jumps, e.g., the 4 lines in Figure 3, because the equivalent noise variance σ^2 increases whenever a datum for \mathbf{X}_S is sampled. (d) For small n , $\text{hhT}(\gamma)$ is closer to FWO, but becomes closer to OV as n increases, as shown by the unmarked solid lines in Figure 3. This is because the theoretical lower bound OV is based on the asymptotically exact single-task OV bound and the $\text{T}(\gamma)$ bound, which is observed to approximate the multi-task learning curve rather closely (point (b)).

Conclusions We have measured the influence of the secondary task on the primary task using the generalization error and the learning curve, parameterizing these with the correlation γ between the two tasks, and the proportion β of observations for the secondary task. We have provided bounds on the generalization error and learning curves, and these bounds highlight the effects of γ and β . This is a step towards understanding the role of the matrix \mathbf{K}_f of inter-task similarities in multi-task GPs with more than two tasks. Analysis on the degenerate case of no training data for the primary task has uncovered an intrinsic limitation of multi-task GP. Our work contributes to an understanding of multi-task learning that is orthogonal to the existing PAC-based results in the literature.

Acknowledgments I thank E Bonilla for motivating this problem, CKI Williams for helpful discussions and for proposing the equivalent isotropic noise approach, and DSO National Laboratories, Singapore, for financial support. This work is supported in part by the EU through the PASCAL2 Network of Excellence.

2 References

- [1] Carl E. Rasmussen and Christopher K. I. Williams. Gaussian Processes for Machine Learning. MIT Press, Cambridge, Massachusetts, 2006.
- [2] Yee Whye Teh, Matthias Seeger, and Michael I. Jordan. Semiparametric latent factor models. In Robert G. Cowell and Zoubin Ghahramani, editors, Proceedings of the 10th International Workshop on Artificial Intelligence and Statistics, pages 333–340. Society for Artificial Intelligence and Statistics, January 2005.
- [3] Edwin V. Bonilla, Felix V. Agakov, and Christopher K. I. Williams. Kernel Multi-task Learning using Task-specific Features. In Marina Meila and Xiaotong Shen, editors, Proceedings of the 11th International Conference on Artificial Intelligence and Statistics. Omni Press, March 2007.
- [4] Kai Yu, Wei Chu, Shipeng Yu, Volker Tresp, and Zhao Xu. Stochastic Relational Models for Discriminative Link Prediction. In B. Schölkopf, J. Platt, and T. Hofmann, editors, Advances in Neural Information Processing Systems 19, Cambridge, MA, 2007. MIT Press.
- [5] Edwin V. Bonilla, Kian Ming A. Chai, and Christopher K.I. Williams. Multi-task Gaussian process prediction. In J.C. Platt, D. Koller, Y. Singer, and S. Roweis, editors, Advances in Neural Information Processing Systems 20. MIT Press, Cambridge, MA, 2008.
- [6] Jonathan Bax-

ter. A Model of Inductive Bias Learning. *Journal of Artificial Intelligence Research*, 12:149?198, March 2000. [7] Andreas Maurer. Bounds for linear multi-task learning. *Journal of Machine Learning Research*, 7:117?139, January 2006. [8] Shai Ben-David and Reba Schuller Borbely. A notion of task relatedness yielding provable multiple-task learning guarantees. *Machine Learning*, 73(3):273?287, 2008. [9] Christopher K. I. Williams and Francesco Vivarelli. Upper and lower bounds on the learning curve for Gaussian processes. *Machine Learning*, 40(1):77?102, 2000. [10] Ya Xue, Xuejun Liao, Lawrence Carin, and Balaji Krishnapuram. Multi-task learning for classification with Dirichlet process prior. *Journal of Machine Learning Research*, 8:35?63, January 2007. [11] Klaus Ritter. Average-Case Analysis of Numerical Problems, volume 1733 of *Lecture Notes in Mathematics*. Springer, 2000. [12] Christopher T. H. Baker. *The Numerical Treatment of Integral Equations*. Clarendon Press, 1977. [13] Peter Sollich and Anason Halees. Learning curves for Gaussian process regression: Approximations and bounds. *Neural Computation*, 14(6):1393?1428, 2002. [14] Noel A. Cressie. *Statistics for Spatial Data*. Wiley, New York, 1993. [15] Manfred Opper and Francesco Vivarelli. General bounds on Bayes errors for regression with Gaussian processes. In Kearns et al. [18], pages 302?308. [16] Giancarlo Ferrari Trecate, Christopher K. I. Williams, and Manfred Opper. Finite-dimensional approximation of Gaussian processes. In Kearns et al. [18], pages 218?224. [17] Huaiyu Zhu, Christopher K. I. Williams, Richard Rohwer, and Michal Morciniec. Gaussian regression and optimal finite dimensional linear models. In Christopher M. Bishop, editor, *Neural Networks and Machine Learning*, volume 168 of *NATO ASI Series F: Computer and Systems Sciences*, pages 167?184. Springer-Verlag, Berlin, 1998. [18] Michael J. Kearns, Sara A. Solla, and David A. Cohn, editors. *Advances in Neural Information Processing Systems 11*, 1999. The MIT Press.