

Documentation: Adaptive Transfer Learning with Gaussian Processes (AT-GP)

Authors of the paper: Bin Cao, Jialin Pan, Yu Zhang, Dit-Yan Yeung, Qiang Yang
Hong Kong University of Science and Technology

Implementation by: Harshdeep J.
SE21UARI044

May 31, 2025

Abstract

Many Transfer Learning methods assume that the source and the target tasks are related, even though there are many tasks which might not be related in reality. In such a case, the knowledge extracted from the source task may not help, or might even hurt the performance of the target task. Thus, to avoid this negative transfer, the authors of the paper propose an Adaptive Transfer Learning approach based on Gaussian Processes which can automatically estimate the similarity between the source task and the target task by proposing a semi-parametric transfer kernel for transfer learning.

The code for this implementation is available on GitHub: [ATGP-implementation](#).

Contents

1	Prerequisite	2
1.1	Gaussian Processes	2
2	Proposed Work	3
3	Implementation	4
3.1	Standard Gaussian Process Implementation	4
3.2	ATGP code implementation	5
4	Experimental Setup and Results	7
4.1	Datasets	7
4.2	Results Summary	7
4.3	Comparison with Paper's Results	9
5	Conclusion	9
6	Acknowledgments	9

1 Prerequisite

1.1 Gaussian Processes

Gaussian Processes are non-parametric Bayesian approach used for regression and probabilistic modeling. This allows us to define a distribution over functions, rather than just parameters, so what GP gives is a distribution of all possible functions that could probably explain the data while giving a measure of uncertainty for its prediction.

1.1.1 Math Behind It

For simplicity, we can assume that the mean function $m(x) = 0$ and the covariance function (kernel) is $k(x, x')$.

A Gaussian process is defined by its mean function and its covariance kernel.

$$f(x) \sim \mathcal{GP}(m(x), k(x, x'))$$

The kernel function measures the similarity between any two input points x and x' . The choice of kernel is crucial as it encodes our assumption about the smoothness, periodicity, or other properties of the function we are trying to model.

A very common kernel is Squared Exponential kernel:

$$k(x_i, x_j) = \sigma_f^2 \exp\left(-\frac{\|x_i - x_j\|^2}{2l^2}\right)$$

Where:

1. x_i and x_j are two input points.
2. $\|x_i - x_j\|^2$ is the squared Euclidean distance between x_i and x_j .
3. σ_f^2 (signal variance) controls the overall vertical scale of the function.
4. l (length scale) controls how quickly the correlation decays with distance. A large l leads to a smoother function.

Now, given a set of training input points $X = [x_1, x_2, \dots, x_N]^T$, the corresponding function values $f(X) = [f(x_1), f(x_2), \dots, f(x_N)]^T$ are jointly Gaussian distributed:

$$f(X) \sim N(\mathbf{0}, K(X, X))$$

Where:

1. $\mathbf{0}$ is a vector of zeros.
2. $K(X, X)$ is a $N \times N$ covariance matrix, where each element $K_{ij} = k(x_i, x_j)$. This matrix captures the prior relationships between all pairs of training data points based on the chosen kernel.

In real-world data, observations are usually noisy. We model this by adding Gaussian noise to the true function values:

$$y_i = f(x_i) + \epsilon$$

Where $\epsilon \sim N(0, \sigma_n^2)$ is independent and identically distributed Gaussian noise with σ_n^2 (noise variance).

So, our observed training outputs y_{train} are also jointly Gaussian distributed:

$$y_{train} \sim N(\mathbf{0}, K(X_{train}, X_{train}) + \sigma_n^2 \mathbf{I})$$

where \mathbf{I} is the identity matrix.

Now, when we get a new test input point x_* , we want to predict its corresponding output y_* . We form a joint distribution of the training outputs y_{train} and the test output $f(x_*)$

$$\begin{bmatrix} y_{train} \\ f(x_*) \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} K(X_{train}, X_{train}) + \sigma_n^2 \mathbf{I} & K(X_{train}, x_*) \\ K(x_*, X_{train}) & k(x_*, x_*) \end{bmatrix} \right)$$

Here, $K(X_{train}, X_{train}) + \sigma_n^2 \mathbf{I}$ is the covariance of the noisy training outputs y_{train} . $K(X_{train}, x_*)$ is the covariance between training points and the test point function value $f(x_*)$, $K(x_*, X_{train})$ is its transpose, and $k(x_*, x_*)$ is the prior variance of the test point function value $f(x_*)$.

While prediction we are interested in the conditional distribution $f(x_*)$, given the observed training data y_{train} . For multivariate Gaussian distribution, the conditional distribution of one block of variables given another is also a Gaussian distribution.

The posterior predictive distribution for $f(x_*)$ is:

$$f(x_*) | X_{train}, y_{train}, x_* \sim N(\mu_{post}, \Sigma_{post})$$

Where:

1. $\mu_{post} = K(x_*, X_{train})(K(X_{train}, X_{train}) + \sigma_n^2 \mathbf{I})^{-1} y_{train}$
2. $\Sigma_{post} = k(x_*, x_*) - K(x_*, X_{train})(K(X_{train}, X_{train}) + \sigma_n^2 \mathbf{I})^{-1} K(X_{train}, x_*)$

Note: $K_{*,train}$ is shorthand for $K(x_*, X_{train})$, $K_{train,train}$ for $K(X_{train}, X_{train})$, and $K_{*,*}$ for $k(x_*, x_*)$.

And while predicting $y_* = f(x_*) + \epsilon_*$, where $\epsilon_* \sim N(0, \sigma_n^2)$, the predictive variance for y_* becomes:

$$\text{Var}(y_* | X_{train}, y_{train}, x_*) = \Sigma_{post} + \sigma_n^2$$

2 Proposed Work

The following modifications are done in standard Gaussian distribution in this work:

1. Instead of modeling the joint distribution of training outputs and test function values and then conditioning on the training outputs y_{train} to obtain the posterior distribution of the test prediction $f(x_*)$, the authors optimize the conditional distribution:

$$p(y^{(T)} | y^{(S)}, X^{(T)}, X^{(S)})$$

2. Instead of using the fixed kernel $K_{ij} = k(x_i, x_j)$, a semi-parametric kernel is introduced:

$$K_{nm} \sim k(x_n, x_m)(2e^{-\zeta(x_n, x_m)\rho} - 1)$$

where $\zeta(x_n, x_m) = 0$ if the points are from the same task, and 1 if not. This adds a task-similarity factor that modulates how much the GP kernel links points from different tasks based on an inferred dissimilarity parameter ρ . Here ρ is modeled as a random variable,

$$\rho \sim \Gamma(b, \mu)$$

where $\Gamma(b, \mu)$ is a Gamma distribution with shape b and mean parameter μ . This gives a closed-form kernel expression:

$$K_{nm} = \begin{cases} k(x_n, x_m) \cdot \left(2 \left(\frac{1}{1+\mu} \right)^b - 1 \right), & \text{if tasks differ} \\ k(x_n, x_m), & \text{if same task} \end{cases}$$

3. The final predictive mean is:

$$m(x) = \sum_{x_j \in X^{(T)}} \alpha_j k(x, x_j) + \sum_{x_i \in X^{(S)}} \lambda \alpha_i k(x, x_i)$$

Where $\lambda = 2(1/(1 + \mu))^b - 1$, quantifies:

- No Transfer when $\lambda \rightarrow 0$,
- Partial Transfer when $0 < |\lambda| < 1$,
- Full Transfer when $\lambda = 1$.

Hence, all the hyperparameters, including μ , b and kernel parameters are learnt by maximizing the conditional likelihood of the target outputs given the source data:

$$\log p(y^{(T)} | y^{(S)}, X^{(T)}, X^{(S)})$$

instead of maximizing the marginal likelihood over all observed data.

3 Implementation

3.1 Standard Gaussian Process Implementation

Kernel Function : SquaredExponentialKernel defines the covariance between any two data points $k(x_i, x_j)$.

```

1 class SquaredExponentialKernel:
2     def __init__(self, length_scale=1.0, sigma_f=1.0):
3         self.length_scale = length_scale
4         self.sigma_f = sigma_f
5
6     def get_params(self):
7         return np.array([self.length_scale, self.sigma_f])
8
9     def set_params(self, params):
10        self.length_scale = params[0]
11        self.sigma_f = params[1]
12
13    def __call__(self, X1, X2):
14        if X1.ndim == 1: X1 = X1[:, np.newaxis]
15        if X2.ndim == 1: X2 = X2[:, np.newaxis]
16
17        if X1.shape[1] != X2.shape[1]:
18            raise ValueError(f"X1 and X2 must have the same number of
19            features. Got {X1.shape[1]} and {X2.shape[1]}")
20
21        sqdist = np.sum(X1**2, 1).reshape(-1, 1) + np.sum(X2**2, 1) - 2 *
22            np.dot(X1, X2.T)
23        sqdist = np.clip(sqdist, 0, np.inf) # Ensure non-negative for sqrt
24            if used elsewhere
25
26        return self.sigma_f**2 * np.exp(-0.5 / self.length_scale**2 *
27            sqdist)

```

Listing 1: SquaredExponentialKernel from atgp.py

Covariance Matrix Calculation : This function applies given kernel to pairs of data points.

```

1 def calculate_covariance_matrix(X1, X2, kernel_func):
2     return kernel_func(X1, X2)

```

Listing 2: Covariance matrix calculation snippet from atgp.py

General GP prediction : This function in SimpleGPR class implements it in the following way.

```

1 # In SimpleGPR class:
2 def predict(self, X_star):
3     # ...
4     K_star = self.kernel(X_star, self.X_train)
5     mean_star = K_star @ self.alpha_
6
7     K_star_star_diag = np.diag(self.kernel(X_star, X_star))
8     # ... v = L^-1 @ K_star.T ...
9     # ... var_reduction_diag = np.sum(v**2, axis=0) ...
10    variance_star_diag = K_star_star_diag + (self._noise_std**2) -
11    var_reduction_diag
12    # ...
13    return mean_star, variance_star_diag.reshape(-1,1)

```

Log Marginal Likelihood for Standard GP : This is the standard objective function maximized for hyperparameter tuning in a regular GP.

```

1 # In SimpleGPR class:
2 def log_marginal_likelihood(self, params_array_opt, X, y):
3     # ... (setup kernel, noise) ...
4     K = calculate_covariance_matrix(X, X, temp_kernel)
5     K_noisy = K + current_noise_var * np.eye(N) + jitter * np.eye(N)
6     # ... (Cholesky, solve for alpha_solve, log_det_K_noisy) ...
7     lml = -0.5 * y.T @ alpha_solve - 0.5 * log_det_K_noisy - N/2.0 * np.log
8     (2 * np.pi)
9     # ...
10    return lml.item()

```

3.2 ATGP code implementation

Conditional Distribution Optimization : The parameters which include kernel hyperparameters, b , μ , and noise terms are optimized by minimizing the negative of this log-likelihood in the ATGP.fit method.

```

1 # In ATGP class:
2 def log_marginal_likelihood_conditional(self, params_array_opt):
3     # ... (parameter setup) ...
4     current_length_scale = params_array_opt[0]
5     current_sigma_f = params_array_opt[1]
6     # ...
7     current_b = params_array_opt[2]
8     current_mu = params_array_opt[3]
9     temp_lambda = self._get_lambda(b_param=current_b, mu_param=current_mu)
10    # ...
11    current_noise_S_var = params_array_opt[4]**2
12    current_noise_T_var = params_array_opt[5]**2
13
14    K11 = calculate_covariance_matrix(self.X_S, self.X_S, temp_kernel)
15    K22 = calculate_covariance_matrix(self.X_T, self.X_T, temp_kernel)
16    K_TS_base = calculate_covariance_matrix(self.X_T, self.X_S, temp_kernel)
17
18    K21 = temp_lambda * K_TS_base
19    K12 = K21.T
20
21    # ... calculation of K11_noisy ...
22    # mu_t = K21 @ inv(K11_noisy) @ y_S

```

```

23     K11_noisy_inv_yS = scipy.linalg.solve_triangular(...)
24     mu_t = K21 @ K11_noisy_inv_yS
25
26     K11_noisy_inv_K12 = scipy.linalg.solve_triangular(...)
27     C_t_main_term = K22 - K21 @ K11_noisy_inv_K12
28     C_t = C_t_main_term + current_noise_T_var * np.eye(K22.shape[0])
29
30     # ...
31     log_likelihood = -0.5 * log_det_Ct - 0.5 * term2_quadratic - len(self.
32     y_T)/2.0 * np.log(2 * np.pi)
33     return log_likelihood.item()
34
35 def fit(self, method='L-BFGS-B', disp=False, maxiter=200):
36     # ...
37     objective = lambda params_opt: -self.
38     log_marginal_likelihood_conditional(params_opt)
39     result = minimize(objective, initial_params, method=method, bounds=
40     bounds, ...)
41     # ...

```

Semi-parametric Kernel & Gamma Distributed Rho : The value of λ is calculated and b and μ are optimized during *fit*.

```

1     # In ATGP class:
2 def __init__(self, ..., initial_b=1.0, initial_mu=1.0, ...):
3     # ...
4     self._b = initial_b
5     self._mu = initial_mu
6     # ...
7
8 def _get_lambda(self, b_param=None, mu_param=None):
9     b_to_use = b_param if b_param is not None else self._b
10    mu_to_use = mu_param if mu_param is not None else self._mu
11    # ... (parameter constraints) ...
12    term_base = 1.0 / (1.0 + mu_to_use)
13    # ... (error handling for pow_term) ...
14    try:
15        pow_term = np.power(term_base, b_to_use)
16        lambda_val = 2 * pow_term - 1
17    # ... (fallback) ...
18    return np.clip(lambda_val, -1.0, 1.0)
19
20 # Used in log_marginal_likelihood_conditional:
21 # K21 = temp_lambda * K_TS_base
22
23 # Used in predict:
24 # K_tilde[:N_S, N_S:] = lambda_val * K_ST_base
25 # K_tilde[N_S:, :N_S] = lambda_val * K_ST_base.T
26 # k_x_star_rows = np.hstack((lambda_val * k_star_S_base, k_star_T_base))
27

```

Predict : This function calculates the final mean.

```

1 # In ATGP class:
2 def predict(self, X_star_T):
3     # ...
4     lambda_val = self._get_lambda()
5     # ...
6     y_all = np.vstack((self.y_S, self.y_T))
7
8     # Construct K_tilde (combined covariance matrix for S and T training
9     data)

```

```

9     K_SS = calculate_covariance_matrix(self.X_S, self.X_S, self.base_kernel
10 )
11     K_TT = calculate_covariance_matrix(self.X_T, self.X_T, self.base_kernel
12 )
13     K_ST_base = calculate_covariance_matrix(self.X_S, self.X_T, self.
14 base_kernel)
15
16     K_tilde = np.zeros((N_all, N_all))
17     K_tilde[:N_S, :N_S] = K_SS
18     K_tilde[N_S:, N_S:] = K_TT
19     K_tilde[:N_S, N_S:] = lambda_val * K_ST_base
20     K_tilde[N_S:, :N_S] = lambda_val * K_ST_base.T
21
22     Big_Lambda_diag = np.concatenate([np.full(N_S, noise_S_var_val),
23                                       np.full(N_T, noise_T_var_val)])
24     Big_Lambda = np.diag(Big_Lambda_diag)
25
26     C_tilde = K_tilde + Big_Lambda + jitter * np.eye(N_all)
27     # ... (alpha calculation) ...
28
29     k_star_S_base = calculate_covariance_matrix(X_star_T, self.X_S, self.
30 base_kernel)
31     k_star_T_base = calculate_covariance_matrix(X_star_T, self.X_T, self.
32 base_kernel)
33     # Apply lambda to the source part of k_x
34     k_x_star_rows = np.hstack((lambda_val * k_star_S_base, k_star_T_base))
35
36     mean_star = k_x_star_rows @ alpha
37     # ... (variance calculation) ...
38     return mean_star, variance_star_diag.reshape(-1,1)
39

```

4 Experimental Setup and Results

Experiments are conducted on three datasets to evaluate the AT-GP model against baseline approaches.

4.1 Datasets

1. Wine Quality:

- (a) Source Task: Predicting quality of white wine (4898 samples).
- (b) Target Task: Predicting quality of red wine (1599 total samples, 5% for training).

2. WiFi:

- (a) Source Task: 500 samples
- (b) Target Task: 500 total samples, 5% for training.

3. SARCOS Robot Arm Inverse Dynamics:

- (a) Source Task: Predicting torque for Joint 1 (subsampling to 1000 samples).
- (b) Target Task: Predicting torque for Joint 2 (subsampling to 1000 total samples, 5% training)

4.2 Results Summary

The following table summarizes the Normalized Mean Squared Error (NMSE) achieved by this implementation. Lower is better.

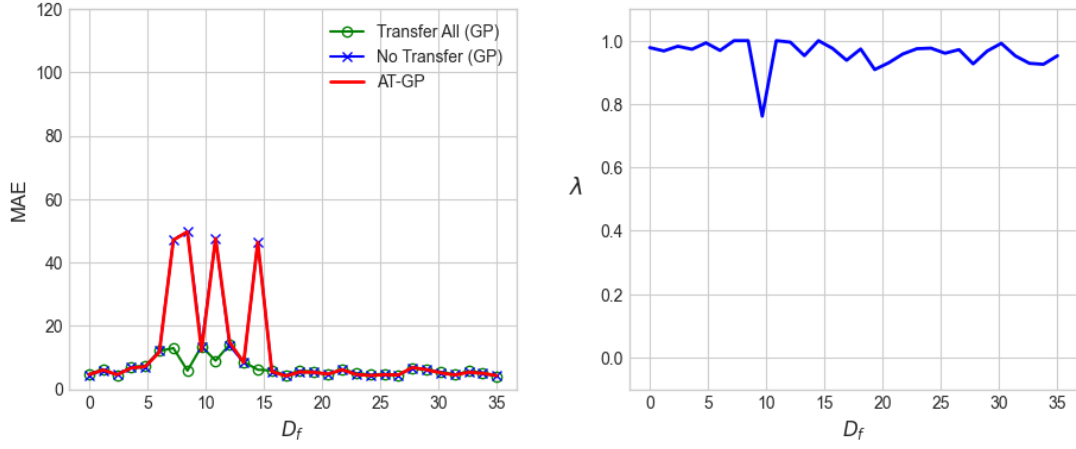


Figure 1: The left figure shows the change to MAE with increasing distance with f . The results are compared with transfer all and no transfer; The right figure shows the change to λ with increasing distance with f . (Own Implementation)

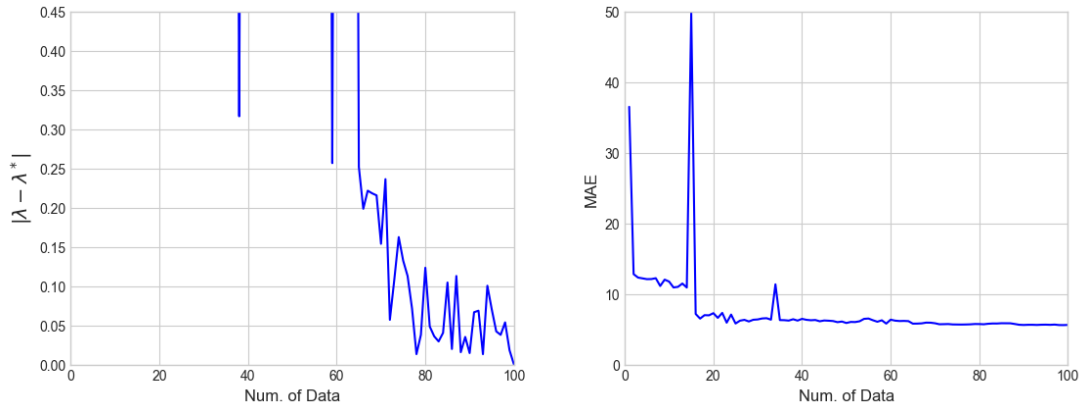


Figure 2: Learning with different numbers of labeled data in the target task. The left figure shows the convergence curve of λ with respect to the number of data. The right figure shows the change to MAE on test data. (Own Implementation)

Table 1: Results from this Implementation

Dataset	No Transfer GP	Transfer All GP	AT-GP
Wine	0.7379 (NMSE)	1.0631 (NMSE)	0.7345 (NMSE)
SARCOS	1.0206 (NMSE)	4.4362 (NMSE)	1.0150 (NMSE)
WiFi	13.3676 (Err Dist)	212.9710 (Err Dist)	13.2415 (Err Dist)

Note: Error Distance(in meters) is the metric for WiFi dataset whereas NMSE is the metric considered for Wine and SARCOS.

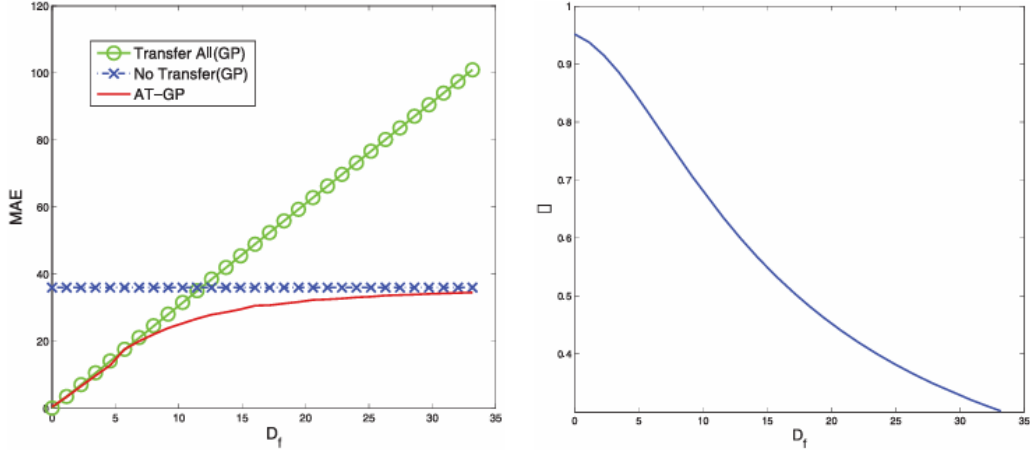


Figure 3: The left figure shows the change to MAE with increasing distance with f . The results are compared with transfer all and no transfer; The right figure shows the change to λ with increasing distance with f . (Original Implementation)

4.3 Comparison with Paper’s Results

The paper (Cao et al., 2010, Table 1) reports the following mean results:

Table 2: Results from Cao et al. (2010, Table 1)

Data	Paper’s No Transfer	Paper’s Transfer All	Paper’s AT-GP
Wine	1.33 (NMSE)	1.37 (NMSE)	1.16 (NMSE)
SARCOS	0.21 (NMSE)	1.58 (NMSE)	0.18 (NMSE)
WiFi	9.18 (Err Dist)	5.28 (Err Dist)	4.98 (Err Dist)

5 Conclusion

According to the method proposed in this paper, how much to transfer is based on how similar the tasks are and negative transfer can be avoided. The experiments on both synthetic and real-world datasets verify the effectiveness of the proposed model **relatively**.

6 Acknowledgments

I acknowledge the assistance of AI language models (ChatGPT and Gemini) during the preparation of this document and the accompanying code. The AI was utilized for tasks such as:

- **Refining**(Not generating from scratch) Python code for the AT-GP implementation, GP implementation is done by self.

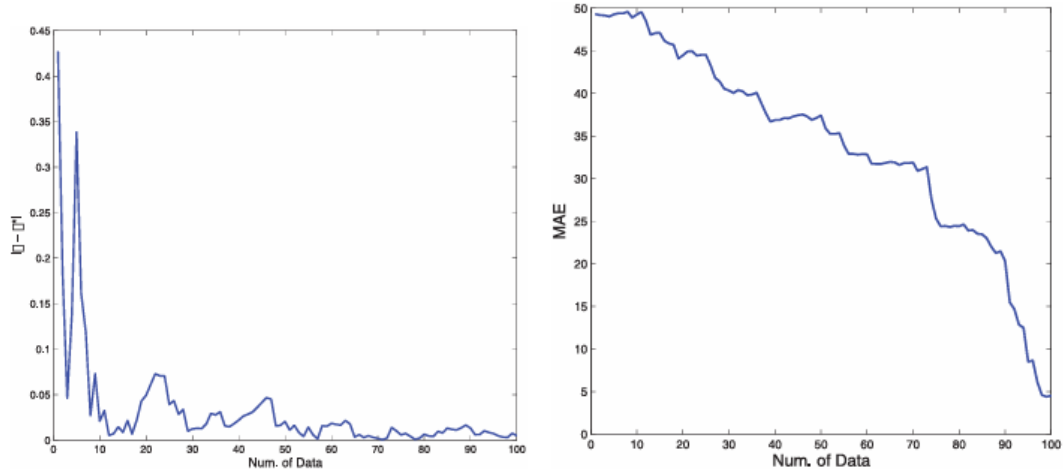


Figure 4: Learning with different numbers of labeled data in the target task. The left figure shows the convergence curve of λ with respect to the number of data. The right figure shows the change to MAE on test data. (Original Implementation)

- Modifying the code to generate the plots shown in this documentation (Code(at-gpr-AI-plotting.ipynb) is separately added in the zip folder).
- Assisting with LaTeX formatting and troubleshooting.