

-Supplementary Material-

Labeling examples that matter: Relevance-Based Active Learning with Gaussian Processes

Alexander Freytag¹, Erik Rodner^{1,2}, Paul Bodesheim¹, and Joachim Denzler¹

¹Computer Vision Group, Friedrich Schiller University Jena, Germany
<http://www.inf-cv.uni-jena.de>

²UC Berkeley ICSI & EECS, United States

Abstract. The following document gives additional details for the paper *Labeling examples that matter: Relevance-Based Active Learning with Gaussian Processes*. In particular, a detailed derivation of the weight vector $\bar{\alpha}$ for Gaussian Process Regression after adding a new example is presented, learning curves of the conducted experiments are visualized, and a qualitative evaluation is given. The information given in this document is not necessary to understand the main paper.

S1 Update of Weight Vector α for GP Regression

Theorem 1 (Closed form update of GP regression weights α).

Let y_* be the class label of a new example \mathbf{x}^* . Let further $\mathbf{K} = \kappa(\mathbf{X}, \mathbf{X})$ be the kernel matrix of all training examples seen so far, $\mathbf{k}_* = \kappa(\mathbf{X}, \mathbf{x}^*)$ is the vector of kernel values between the new example and all training examples, and $k_{**} = \kappa(\mathbf{x}^*, \mathbf{x}^*)$ is the self-similarity value of the new sample. In addition, we denote with \mathbf{y} the vector of labels for all training samples, and σ_n^2 indicates the noise parameter for model regularization. Finally, let $\alpha = \mathbf{K}^{-1}\mathbf{y}$ be the weight vector for GP regression. Then we can compute the weight vector $\bar{\alpha}$ after adding \mathbf{x}^* to the training set as follows:

$$\begin{aligned} \bar{\alpha} &= \bar{\mathbf{K}}^{-1} \begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} = \begin{bmatrix} \alpha \\ 0 \end{bmatrix} + \underbrace{\frac{1}{\sigma_{f_*}^2 + \sigma_n^2}}_{\langle 1 \rangle} \underbrace{\begin{bmatrix} (\mathbf{K} + \sigma_n^2 \cdot \mathbf{I})^{-1} \mathbf{k}_* \\ -1 \end{bmatrix}}_{\langle 2 \rangle} \underbrace{(\mathbf{k}_*^T \alpha - y_*)}_{\langle 3 \rangle} \quad (\text{T1}) \\ &= \begin{bmatrix} \alpha \\ 0 \end{bmatrix} + \frac{\mu_* - y_*}{\sigma_{f_*}^2 + \sigma_n^2} \begin{pmatrix} \mathbf{K}^{-1} \mathbf{k}_* \\ -1 \end{pmatrix}. \end{aligned}$$

The three factors can be interpreted as described in the main paper.

Proof.

$$\begin{aligned}
\bar{\alpha} &= \bar{\mathbf{K}}^{-1} \begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \\
&= \begin{bmatrix} \mathbf{K} & \mathbf{k}_* \\ \mathbf{k}_*^T & k_{**} + \sigma_n^2 \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \\
&\stackrel{[2, \text{Eq. (2.8.18), p. 117}]}{=} \begin{bmatrix} \left(\mathbf{K} - \frac{1}{k_{**} + \sigma_n^2} \mathbf{k}_* \mathbf{k}_*^T \right)^{-1} & - \left(\mathbf{K} - \frac{1}{k_{**} + \sigma_n^2} \mathbf{k}_* \mathbf{k}_*^T \right)^{-1} \frac{1}{k_{**} + \sigma_n^2} \mathbf{k}_* \\ - \left(k_{**} + \sigma_n^2 - \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{k}_* \right)^{-1} \mathbf{k}_*^T \mathbf{K}^{-1} & \left(k_{**} + \sigma_n^2 - \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{k}_* \right)^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \\
&\stackrel{[2, \text{Eq. (2.8.19), p. 117}]}{=} \begin{bmatrix} \mathbf{K}^{-1} + \mathbf{K}^{-1} \mathbf{k}_* \left(k_{**} + \sigma_n^2 - \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{k}_* \right)^{-1} \mathbf{k}_*^T \mathbf{K}^{-1} & - \left(\mathbf{K}^{-1} + \mathbf{K}^{-1} \mathbf{k}_* \left(k_{**} + \sigma_n^2 - \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{k}_* \right)^{-1} \mathbf{k}_*^T \mathbf{K}^{-1} \right) \frac{1}{k_{**} + \sigma_n^2} \mathbf{k}_* \\ - \left(k_{**} + \sigma_n^2 - \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{k}_* \right)^{-1} \mathbf{k}_*^T \mathbf{K}^{-1} & \left(k_{**} + \sigma_n^2 - \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{k}_* \right)^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \\
&\stackrel{\text{Definition of } \alpha \text{ and GP variance}}{=} \begin{pmatrix} \alpha + \mathbf{K}^{-1} \mathbf{k}_* \frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \mathbf{k}_*^T \alpha - \left(\mathbf{K}^{-1} + \mathbf{K}^{-1} \mathbf{k}_* \frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \mathbf{k}_*^T \mathbf{K}^{-1} \right) \frac{y_*}{k_{**} + \sigma_n^2} \mathbf{k}_* \\ - \frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \mathbf{k}_*^T \alpha + \frac{y_*}{\sigma_{f_*}^2 + \sigma_n^2} \end{pmatrix} \\
&\stackrel{\text{Factor out } \alpha}{=} \begin{bmatrix} \alpha \\ 0 \end{bmatrix} + \begin{pmatrix} \mathbf{K}^{-1} \mathbf{k}_* \frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \mathbf{k}_*^T \alpha - \frac{y_*}{k_{**} + \sigma_n^2} \left(\mathbf{K}^{-1} + \mathbf{K}^{-1} \mathbf{k}_* \frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \mathbf{k}_*^T \mathbf{K}^{-1} \right) \mathbf{k}_* \\ - \frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \mathbf{k}_*^T \alpha + \frac{y_*}{\sigma_{f_*}^2 + \sigma_n^2} \end{pmatrix} \\
&\stackrel{\text{Simplify denominator}}{=} \begin{bmatrix} \alpha \\ 0 \end{bmatrix} + \begin{pmatrix} \mathbf{K}^{-1} \mathbf{k}_* \frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \mathbf{k}_*^T \alpha - \frac{y_*}{k_{**} + \sigma_n^2} \left(\mathbf{K}^{-1} + \mathbf{K}^{-1} \mathbf{k}_* \frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \mathbf{k}_*^T \mathbf{K}^{-1} \right) \mathbf{k}_* \\ \frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \left(y_* - \mathbf{k}_*^T \alpha \right) \end{pmatrix} \\
&\stackrel{\text{Expand numerator with } \mathbf{k}_*}{=} \begin{bmatrix} \alpha \\ 0 \end{bmatrix} + \begin{pmatrix} \mathbf{K}^{-1} \mathbf{k}_* \frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \mathbf{k}_*^T \alpha - \frac{y_*}{k_{**} + \sigma_n^2} \left(\mathbf{K}^{-1} \mathbf{k}_* + \mathbf{K}^{-1} \mathbf{k}_* \frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{k}_* \right) \\ \frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \left(y_* - \mathbf{k}_*^T \alpha \right) \end{pmatrix} \\
&= (*)
\end{aligned}$$

$$\begin{aligned}
(*) &= \begin{bmatrix} \boldsymbol{\alpha} \\ 0 \end{bmatrix} + \left(\mathbf{K}^{-1} \mathbf{k}_* \frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \mathbf{k}_*^T \boldsymbol{\alpha} - \frac{y_*}{k_{**} + \sigma_n^2} \left(\mathbf{K}^{-1} \mathbf{k}_* + \mathbf{K}^{-1} \mathbf{k}_* \frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{k}_* \right) \right) \\
&\quad \frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \left(y_* - \mathbf{k}_*^T \boldsymbol{\alpha} \right) \\
\text{Factor out } \mathbf{K}^{-1} \mathbf{k}_* \text{ in enumerator} &= \begin{bmatrix} \boldsymbol{\alpha} \\ 0 \end{bmatrix} + \left(\mathbf{K}^{-1} \mathbf{k}_* \left(\frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \mathbf{k}_*^T \boldsymbol{\alpha} - \frac{y_*}{k_{**} + \sigma_n^2} \left(1 + \frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{k}_* \right) \right) \right) \\
&\quad \frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \left(y_* - \mathbf{k}_*^T \boldsymbol{\alpha} \right) \\
\text{Definition of GP variance} &= \begin{bmatrix} \boldsymbol{\alpha} \\ 0 \end{bmatrix} + \left(\mathbf{K}^{-1} \mathbf{k}_* \left(\frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \mathbf{k}_*^T \boldsymbol{\alpha} - \frac{y_*}{k_{**} + \sigma_n^2} \left(1 + \frac{k_{**} - \sigma_{f_*}^2}{\sigma_{f_*}^2 + \sigma_n^2} \right) \right) \right) \\
&\quad \frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \left(y_* - \mathbf{k}_*^T \boldsymbol{\alpha} \right) \\
\text{Expand 1 in enumerator} &= \begin{bmatrix} \boldsymbol{\alpha} \\ 0 \end{bmatrix} + \left(\mathbf{K}^{-1} \mathbf{k}_* \left(\frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \mathbf{k}_*^T \boldsymbol{\alpha} - \frac{y_*}{k_{**} + \sigma_n^2} \left(\frac{\sigma_{f_*}^2 + \sigma_n^2}{\sigma_{f_*}^2 + \sigma_n^2} + \frac{k_{**} - \sigma_{f_*}^2}{\sigma_{f_*}^2 + \sigma_n^2} \right) \right) \right) \\
&\quad \frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \left(y_* - \mathbf{k}_*^T \boldsymbol{\alpha} \right) \\
\text{Merge denominators} &= \begin{bmatrix} \boldsymbol{\alpha} \\ 0 \end{bmatrix} + \left(\mathbf{K}^{-1} \mathbf{k}_* \left(\frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \mathbf{k}_*^T \boldsymbol{\alpha} - \frac{y_*}{k_{**} + \sigma_n^2} \left(\frac{k_{**} + \sigma_n^2}{\sigma_{f_*}^2 + \sigma_n^2} \right) \right) \right) \\
&\quad \frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \left(y_* - \mathbf{k}_*^T \boldsymbol{\alpha} \right) \\
\text{Simplify enumerator} &= \begin{bmatrix} \boldsymbol{\alpha} \\ 0 \end{bmatrix} + \left(\mathbf{K}^{-1} \mathbf{k}_* \left(\frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \mathbf{k}_*^T \boldsymbol{\alpha} - \frac{y_*}{\sigma_{f_*}^2 + \sigma_n^2} \right) \right) \\
&\quad \frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \left(y_* - \mathbf{k}_*^T \boldsymbol{\alpha} \right) \\
\text{Factor out } \frac{1}{\sigma_{f_*}^2 + \sigma_n^2} &= \begin{bmatrix} \boldsymbol{\alpha} \\ 0 \end{bmatrix} + \frac{1}{\sigma_{f_*}^2 + \sigma_n^2} \left(\mathbf{K}^{-1} \mathbf{k}_* \left(\mathbf{k}_*^T \boldsymbol{\alpha} - y_* \right) \right) \\
&\quad \left(y_* - \mathbf{k}_*^T \boldsymbol{\alpha} \right) \\
\text{Factor out } \mathbf{k}_*^T \boldsymbol{\alpha} - y_* &= \begin{bmatrix} \boldsymbol{\alpha} \\ 0 \end{bmatrix} + \frac{\mathbf{k}_*^T \boldsymbol{\alpha} - y_*}{\sigma_{f_*}^2 + \sigma_n^2} \begin{bmatrix} \mathbf{K}^{-1} \mathbf{k}_* \\ -1 \end{bmatrix} \\
\text{Definition of GP mean and variance} &= \begin{bmatrix} \boldsymbol{\alpha} \\ 0 \end{bmatrix} + \frac{\mu_* - y_*}{\sigma_*^2} \begin{bmatrix} \mathbf{K}^{-1} \mathbf{k}_* \\ -1 \end{bmatrix} \quad \square
\end{aligned}$$

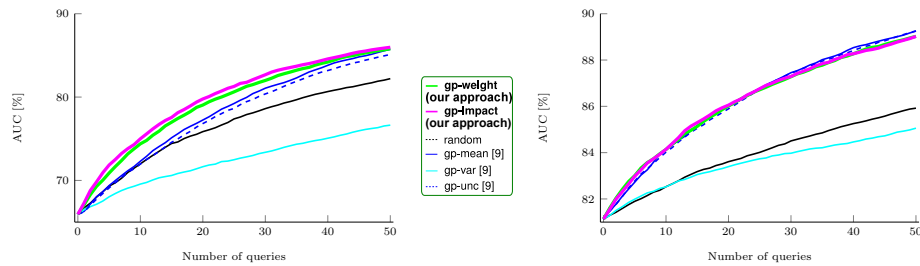


Fig. 1. Active Learning on **Caltech-256**. Initial number of samples per class is 1 and 5, respectively.

S2 Recognition Rates of Experimental Results

In the main paper, we presented active learning gains over random sampling for two popular image categorization datasets. Due to the lack of space, we moved the visualization of learning curves to this supplementary material document.

S2.1 Results on Caltech-256

Active learning curves for Caltech-256 are given in Fig. 1. The maximum accuracy when using all unlabeled data as additional training images is 93.81% AUC. As in the main document, results were obtained by averaging over the corresponding 1,000 binary settings. Hyperparameters were optimized by maximizing the marginal likelihood.

S2.2 Results on ImageNet

Active learning curves for ImageNet (2010 challenge) are given in Fig. 2. The maximum accuracy when using all unlabeled data as additional training images is 89.2% AUC. As in the main document, results were obtained by averaging over the corresponding 1,000 binary settings. Hyperparameters were optimized by maximizing the marginal likelihood.

S3 Queried Images of a Single Run on ImageNet

For a qualitative evaluation, we presented queried images of a single active learning setup in the main document. Due to the lack of space, we only included results of $\mathcal{Q}_{\text{impact}}$ and $\mathcal{Q}_{\sigma_2^*}$, and moved the remaining results to this document. The queried images of all strategies are given in Fig. 3.

References

2. Bernstein, D.S.: Matrix Mathematics. Princeton University Press, 2nd edn. (2009)
9. Kapoor, A., Grauman, K., Urtasun, R., Darrell, T.: Gaussian processes for object categorization. IJCV 88, 169–188 (2010)

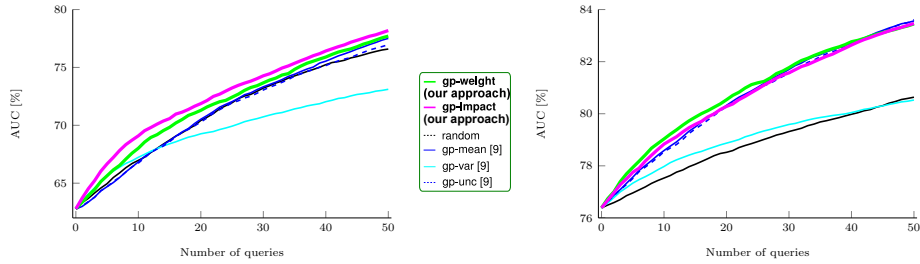


Fig. 2. Active Learning on **ImageNet**. Initial number of samples per class is 1 and 5, respectively.



Fig. 3. Queried images of GP-based active learning methods: (1) random queries, (2) Q_{μ^*} , (3) $Q_{\sigma_2^2}$, and (4) Q_{unc} introduced by [9] as well as our strategies (5) Q_{weight} and (6) Q_{impact} (from top to bottom). Green and blue borders indicate images of positive and negative classes, respectively.