
Correction to “Stability and Hypothesis Transfer Learning”

Ilja Kuzborskij

Idiap Research Institute, Switzerland
 École Polytechnique Fédérale de Lausanne (EPFL), Switzerland

ILJA.KUZBORSKIJ@IDIAP.CH

Francesco Orabona

Toyota Technological Institute at Chicago, USA

FRANCESCO@ORABONA.COM

Abstract

There is an error in “Stability and Hypothesis Transfer Learning” (Kuzborskij & Orabona, 2013) which appeared in proceedings of ICML 2013. The Leave-One-Out generalization bound for Hypothesis Transfer Learning algorithm through Regularized Least Squares with biased regularization does not have the right convergence rate with respect to the regularization parameter λ and the source risk on the target domain, $R_\mu(f')$. This erratum describes the error and proves the correct generalization guarantees. The correct rate is in $\mathcal{O}(\frac{1}{m\lambda^{1.5}})$, instead of incorrectly claimed $\mathcal{O}(\frac{1}{m\lambda})$. However, the correct rate is still better than the usual one for Regularized Least Squares obtained via algorithmic stability analysis. Finally, corrected analysis still preserves the main contribution, that is, the relatedness of the source and target domains accelerates the convergence of the Leave-One-Out error to the generalization error.

1. Description of Error

The error was committed in the proof of Theorem 3 (Kuzborskij & Orabona, 2013), where Lemma 2 was applied incorrectly. This rendered Theorem 2, proving generalization guarantees for Hypothesis Transfer Learning (HTL) algorithm analyzed, invalid. Theorem 3 proves an upper-bound on hypothesis stability (Bousquet & Elisseeff, 2002) with respect to the square loss,

$$\forall i \in \{1, \dots, m\},$$

$$\mathbb{E}_{S,(\mathbf{x},y)} [|(f_S(\mathbf{x}) - y)^2 - (f_{S^{\setminus i}}(\mathbf{x}) - y)^2|] \leq \gamma.$$

Here (\mathbf{x}, y) assumed to be any example drawn i.i.d. from p.d.f. μ . However, Lemma 2 proves an upper

bound on the quantity $(f_{S^{\setminus i}}(\mathbf{x}_i) - y_i)^2$, where $(\mathbf{x}_i, y_i) \in S$, in other words, belongs to the training set. That is, in Theorem 3, $(f_{S^{\setminus i}}(\mathbf{x}_i) - y_i)^2$ appears instead of $(f_{S^{\setminus i}}(\mathbf{x}) - y)^2$.

To use the correct quantity $(f_{S^{\setminus i}}(\mathbf{x}) - y)^2$, we first obtain a closed form solution to $f_{S^{\setminus i}}(\mathbf{x})$, considering linear hypotheses $f(\mathbf{x}) := \mathbf{x}^\top \mathbf{w}$ and $f_{S^{\setminus i}}(\mathbf{x}) := \mathbf{x}^\top \mathbf{w}_{S^{\setminus i}}$. The derivation is given in Lemma 1.

A very similar error also appears in the proof of Lemma 4 (Kuzborskij & Orabona, 2013), first result. The nature of error is the same, and we fix it in the proof of the Theorem 1.

Lemma 1. *Let \mathbf{w}_S be the hypothesis produced by the Regularized Least Squares (RLS) algorithm given training set S . For any sample $(\mathbf{x}, y) \stackrel{i.i.d.}{\sim} \mu$ and $(\mathbf{x}_i, y_i) \in S$, such that $\|\mathbf{x}\|, \|\mathbf{x}_i\| \leq 1$, we have that the hypothesis $\mathbf{w}_{S^{\setminus i}}$ produced by the same RLS algorithm on a training set $S^{\setminus i} \forall i \in \{1, \dots, m\}$, satisfies*

$$|\mathbf{x}^\top \mathbf{w}_S - \mathbf{x}^\top \mathbf{w}_{S^{\setminus i}}| \leq \frac{1}{m\lambda} |\mathbf{x}_i^\top \mathbf{w}_{S^{\setminus i}} - y_i|.$$

Proof. Define $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_m]$, $\mathbf{M} = \mathbf{X}^\top \mathbf{X} + m\lambda \mathbf{I}$. It is straightforward to see that $\mathbf{x}^\top \mathbf{w}_S$ is equal to

$$\begin{bmatrix} \mathbf{x}^\top \mathbf{X} & \mathbf{x}^\top \mathbf{x}_i \end{bmatrix} \begin{bmatrix} \mathbf{M} & \mathbf{X}^\top \mathbf{x}_i \\ \mathbf{x}_i^\top \mathbf{X} & \|\mathbf{x}_i\|^2 + m\lambda \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{y} \\ y_i \end{bmatrix}. \quad (1)$$

Expanding the middle term and using the block-wise matrix inversion property (Petersen & Pedersen, 2008) we get

$$\begin{aligned} \begin{bmatrix} \mathbf{M} & \mathbf{X}^\top \mathbf{x}_i \\ \mathbf{x}_i^\top \mathbf{X} & \|\mathbf{x}_i\|^2 + m\lambda \end{bmatrix}^{-1} &= \begin{bmatrix} \mathbf{M}^{-1} & \mathbf{0} \\ \mathbf{0}^\top & 0 \end{bmatrix} \\ &+ \frac{1}{a} \begin{bmatrix} \mathbf{M}^{-1} \mathbf{X}^\top \mathbf{x}_i \\ -1 \end{bmatrix} \begin{bmatrix} \mathbf{x}_i^\top \mathbf{X} \mathbf{M}^{-1} & -1 \end{bmatrix}, \end{aligned}$$

where $a := \|\mathbf{x}_i\|^2 + m\lambda - \mathbf{x}_i^\top \mathbf{X} \mathbf{M}^{-1} \mathbf{X}^\top \mathbf{x}_i$. Plugging

this result into (1) yields

$$\mathbf{x}^\top \mathbf{w}_S = \mathbf{x}^\top \mathbf{w}_{S^{\setminus i}} + \frac{\mathbf{x}^\top \left(\mathbf{I} - \mathbf{X} \mathbf{M}^{-1} \mathbf{X}^\top \right) \mathbf{x}_i}{a} (y_i - \mathbf{x}_i^\top \mathbf{w}_{S^{\setminus i}}).$$

Using result of Lemma 2, we have that $m\lambda \leq a$ and in addition by Cauchy-Schwarz inequality we have that $\mathbf{x} \left(\mathbf{I} - \mathbf{X} \mathbf{M}^{-1} \mathbf{X}^\top \right) \mathbf{x}_i \leq 1$, since $\|\mathbf{x}\|, \|\mathbf{x}_i\| \leq 1$. \square

Lemma 2. For all $\mathbf{X} \in \mathbb{R}^{m \times d}$, $m, \lambda \geq 0$, we have that the matrix

$$\mathbf{I} - \mathbf{X} \left(\mathbf{X}^\top \mathbf{X} + m\lambda \mathbf{I} \right)^{-1} \mathbf{X}^\top$$

is PSD and its maximum eigenvalue is less than 1.

Proof.

$$\begin{aligned} & \mathbf{I} - \mathbf{X} \left(\mathbf{X}^\top \mathbf{X} + m\lambda \mathbf{I} \right)^{-1} \mathbf{X}^\top \\ &= \mathbf{I} - \left(\mathbf{X} \mathbf{X}^\top + m\lambda \mathbf{I} \right)^{-1} \mathbf{X} \mathbf{X}^\top \quad (2) \\ &= \mathbf{I} - \mathbf{U} \left(\boldsymbol{\Sigma} + m\lambda \mathbf{I} \right)^{-1} \mathbf{U}^\top \mathbf{U} \boldsymbol{\Sigma} \mathbf{U}^\top \\ &= \mathbf{U} \left(\mathbf{I} - \left(\boldsymbol{\Sigma} + m\lambda \mathbf{I} \right)^{-1} \boldsymbol{\Sigma} \right) \mathbf{U}^\top, \end{aligned}$$

where we used identity $\left(\mathbf{X} \mathbf{X}^\top + m\lambda \mathbf{I} \right)^{-1} \mathbf{X} = \mathbf{X} \left(\mathbf{X}^\top \mathbf{X} + m\lambda \mathbf{I} \right)^{-1}$ to obtain (2). Subsequently we use decomposition $\mathbf{X} \mathbf{X}^\top = \mathbf{U} \boldsymbol{\Sigma} \mathbf{U}^\top$. \square

2. Weaker HTL Guarantees through Correct Stability

The following HTL algorithm is analyzed by Kuzborskij & Orabona (2013).

Algorithm 1. RLS transfer algorithm by altering training set as $\{(\mathbf{x}_i, y_i - f'(\mathbf{x}_i)) : 1 \leq i \leq m\}$ produces a hypothesis

$$f_S^{htl'}(\mathbf{x}) = T_C(\mathbf{x}^\top \hat{\mathbf{w}}_S) + f'(\mathbf{x}),$$

where

$$\hat{\mathbf{w}}_S := \underset{\mathbf{u}}{\operatorname{argmin}} \frac{1}{m} \sum_{i=1}^m (\mathbf{u}^\top \mathbf{x}_i - y_i + f'(\mathbf{x}_i))^2 + \lambda \|\mathbf{u}\|^2,$$

and the truncation function $T_C(\hat{y})$ is defined as $T_C(\hat{y}) = \min(\max(\hat{y}, -C), C)$.

We prove Theorem 1 which is corrected version of originally presented Theorem 2 in (Kuzborskij & Orabona, 2013).

Theorem 1. Set $\lambda \geq \frac{1}{m}$. If $C \geq B + \|f'\|_\infty$, then for Algorithm 1 we have

$$\begin{aligned} & \mathbb{E}_S[(R_\mu(f_S^{htl'}) - \hat{R}^{\text{loo}}(f_S^{htl'}))^2] \\ &= \mathcal{O} \left(\frac{C^2 \sqrt{R_\mu(f') T_{C^2} \left(\frac{R_\mu(f')}{\lambda} \right) + R_\mu(f')^2}}{m\lambda^{1.5}} \right). \end{aligned}$$

If $C = \infty$, then for Algorithm 1 we have

$$\begin{aligned} & \mathbb{E}_S[(R_\mu(f_S^{htl'}) - \hat{R}^{\text{loo}}(f_S^{htl'}))^2] \\ &= \mathcal{O} \left(\frac{R_\mu(f') (\|f'\|_\infty + B)^2}{m\lambda^3} \right). \end{aligned}$$

2.1. Implications of Weaker Result

The first and the main difference comes in weaker bound on the second order moment of the Leave-One-Out (LOO) error when one applies the truncation on the predictions in the range $[-B; B]$. In (Kuzborskij & Orabona, 2013) we have claimed it to be in $\mathcal{O}\left(\frac{B^2}{m\lambda}\right)$. However in the following we obtain weaker corrected rate in $\mathcal{O}\left(\frac{B^2}{m\lambda^{1.5}}\right)$. Note that this rate is still better than the one that can be obtained for RLS through stability analysis, $\mathcal{O}\left(\frac{1}{m\lambda^3}\right)$ (De Vito et al., 2005).

The second difference comes in dependence on the risk of the source on the target domain. The correct bound is in $\mathcal{O}\left(\frac{\sqrt{R_\mu(f')^2 + R_\mu(f')}}{m\lambda^{1.5}}\right)$, rather than $\mathcal{O}\left(\frac{R_\mu(f')}{m\lambda}\right)$, incorrectly suggested earlier. Nevertheless, the important behavior of HTL is still present, that is whenever the source hypothesis f' performs well on the target domain, in other words, for $\frac{\sqrt{R_\mu(f')^2 + R_\mu(f')}}{\lambda^{1.5}} \rightarrow 0$, the LOO error approaches expected risk with probability 1.

2.2. Proof of Theorem 1

In this section we give the proof of the main theorem. First we prove utility Lemma 3.

Lemma 3. $\forall a, b, \hat{y} \in \mathbb{R}$,

$$|(a - \hat{y})^2 - (b - \hat{y})^2| \leq (a - b)^2 + 2|(b - \hat{y})(a - b)|.$$

Proof.

$$\begin{aligned} & |(a - \hat{y})^2 - (b - \hat{y})^2| \\ &= |a^2 - b^2 - 2\hat{y}(a - b)| \\ &= |(a - b)^2 - 2b^2 + 2ab - 2\hat{y}(a - b)| \\ &= |(a - b)^2 + 2(b - \hat{y})(a - b)| \\ &\leq (a - b)^2 + 2|(b - \hat{y})(a - b)|. \end{aligned}$$

\square

The following Theorem upper-bounds the hypothesis stability of the Algorithm 1.

Theorem 2. *The hypothesis stability of Algorithm 1 is upper bounded as*

$$\begin{aligned} \gamma &\leq T_{4C^2} \left(\frac{2R_\mu(f')}{m^2\lambda^2} \left(1 + \frac{1}{\lambda} \right) \right) \\ &\quad + 2T_{2C} \left(\frac{\sqrt{2R_\mu(f')}}{m\lambda} \sqrt{1 + \frac{1}{\lambda}} \right) \\ &\quad \cdot \sqrt{2T_{C^2} \left(\frac{R_\mu(f')}{\lambda} \right) + 2R_\mu(f')} . \end{aligned}$$

Proof. From Lemma 3 with $a = T_C(\Delta + \epsilon)$, $b = T_C(\Delta)$, $\hat{y} = y - f'(\mathbf{x})$ and also using the fact that $|T_C(\Delta + \epsilon) - T_C(\Delta)| \leq \min(|\epsilon|, 2C)$, we have

$$\begin{aligned} &|(T_C(\Delta + \epsilon) - y + f'(\mathbf{x}))^2 - (T_C(\Delta) - y + f'(\mathbf{x}))^2| \\ &\leq \min(\epsilon^2, 4C^2) + 2 \min(|\epsilon|, 2C) |T_C(\Delta) - y + f'(\mathbf{x})| . \end{aligned}$$

Set $\Delta := \mathbf{x}^\top \mathbf{w}_{S^{\setminus i}}$, and $\Delta + \epsilon := \mathbf{x}^\top \mathbf{w}_S$. Taking the expectation $\mathbb{E}[\cdot] = \mathbb{E}_{S, (\mathbf{x}, y)}[\cdot]$, and using Jensen’s and Cauchy-Schwarz’s inequalities, we have

$$\begin{aligned} &\mathbb{E} [|(T_C(\Delta + \epsilon) - y + f'(\mathbf{x}))^2 - (T_C(\Delta) - y + f'(\mathbf{x}))^2|] \\ &\leq \min(\mathbb{E}[\epsilon^2], 4C^2) + 2 \min\left(\sqrt{\mathbb{E}[\epsilon^2]}, 2C\right) . \end{aligned}$$

$$\begin{aligned} &\quad \cdot \sqrt{\mathbb{E}[(T_C(\Delta) - y + f'(\mathbf{x}))^2]} \\ &\leq \min(\mathbb{E}[\epsilon^2], 4C^2) + 2 \min\left(\sqrt{\mathbb{E}[\epsilon^2]}, 2C\right) \\ &\quad \cdot \sqrt{\mathbb{E}[2T_{C^2}(\|\hat{\mathbf{w}}_S\|^2) + 2(f'(\mathbf{x}) - y)^2]} \quad (3) \end{aligned}$$

$$\begin{aligned} &\leq \min(\mathbb{E}[\epsilon^2], 4C^2) + 2 \min\left(\sqrt{\mathbb{E}[\epsilon^2]}, 2C\right) \\ &\quad \cdot \sqrt{2T_{C^2} \left(\frac{R_\mu(f')}{\lambda} \right) + 2R_\mu(f')} . \quad (4) \end{aligned}$$

In (3) we apply Cauchy-Schwarz inequality and elementary inequality $(a + b)^2 \leq 2a^2 + 2b^2$, while (4) comes from the first result of Lemma 3 in (Kuzborskij & Orabona, 2013).

We now use the fact that

$$\begin{aligned} \mathbb{E}[\epsilon^2] &\leq \frac{1}{m^2\lambda^2} \mathbb{E}[(\mathbf{x}^\top \hat{\mathbf{w}}_{S^{\setminus i}} - y_i + f'(\mathbf{x}))^2] \\ &\leq \frac{2}{m^2\lambda^2} \mathbb{E}[\|\hat{\mathbf{w}}_{S^{\setminus i}}\|^2 + (y - f'(\mathbf{x}))^2] \\ &\leq \frac{2R_\mu(f')}{m^2\lambda^2} \left(\frac{m-1}{m} \frac{1}{\lambda} + 1 \right) \\ &\leq \frac{2R_\mu(f')}{m^2\lambda^2} \left(\frac{1}{\lambda} + 1 \right) . \end{aligned}$$

Putting all together we have

$$\begin{aligned} &\mathbb{E}_{S, (\mathbf{x}, y)} [|(T_C(\Delta) - y + f'(\mathbf{x}))^2 - (T_C(\Delta + \epsilon) - y + f'(\mathbf{x}))^2|] \\ &\leq T_{4C^2} \left(\frac{2R_\mu(f')}{m^2\lambda^2} \left(1 + \frac{1}{\lambda} \right) \right) \\ &\quad + 2T_{2C} \left(\frac{\sqrt{2R_\mu(f')}}{m\lambda} \sqrt{1 + \frac{1}{\lambda}} \right) \\ &\quad \cdot \sqrt{2T_{C^2} \left(\frac{R_\mu(f')}{\lambda} \right) + 2R_\mu(f')} . \end{aligned}$$

□

Proof of Theorem 1. We apply Theorem 1 from (Kuzborskij & Orabona, 2013). To apply this theorem, we need to upper-bound quantities $M, \mathbb{E}_S[\ell(\mathbf{w}_{S^{\setminus i}}, (\mathbf{x}_i, y_i))]$ and γ . As $\mathbb{E}_S[\ell(\mathbf{w}_{S^{\setminus i}}, (\mathbf{x}_i, y_i))]$ is already correctly bounded in (Kuzborskij & Orabona, 2013) as

$$\begin{aligned} &\mathbb{E}_S[\ell(f_{S^{\setminus i}}, (\mathbf{x}_i, y_i))] \\ &\leq 2 \left(1 + \frac{1}{m\lambda} \right)^2 \left(T_{C^2} \left(\frac{R_\mu(f')}{\lambda} \right) + R_\mu(f') \right) , \end{aligned}$$

we use bound on γ given by Theorem 2.

By definition of Theorem 1 in (Kuzborskij & Orabona, 2013),

$$\ell(\mathbf{w}_{S^{\setminus i}}, (\mathbf{x}, \hat{y})) \leq M, \quad \forall \mathbf{x} \in \mathcal{X}, \hat{y} \in \mathcal{Y} .$$

So we have

$$\begin{aligned} &\sup_{\mathbf{x}, y} (T_C(\mathbf{x}^\top \mathbf{w}_{S^{\setminus i}}) - \hat{y} + f'(\mathbf{x}))^2 \\ &\leq \left(T_C \left(\frac{B + \|f'\|_\infty}{\sqrt{\lambda}} \right) + B + \|f'\|_\infty \right)^2 . \end{aligned}$$

We have this result, because term $T_C(\mathbf{x}^\top \mathbf{w}_{S^{\setminus i}})$ can be simultaneously upper-bounded by C and, using Cauchy-Schwarz inequality, $\|\mathbf{w}_{S^{\setminus i}}\|$. Consequently, $\|\mathbf{w}_{S^{\setminus i}}\|$ is bounded using second result of Lemma 3 in (Kuzborskij & Orabona, 2013).

Putting it all together and applying Theorem 1 (Kuzborskij & Orabona, 2013), we have the stated result. The dominant rates in $\mathcal{O}(\cdot)$ notation, in both truncated and untruncated cases, come from the bound on the component $M\gamma$ in Theorem 1 in (Kuzborskij & Orabona, 2013). \square

3. Conclusions

For an Algorithm 1 we have obtained the generalization bound of the form,

$$R_\mu(f_S^{htt'}) \leq \hat{R}^{\text{loo}}(f_S^{htt'}) + \mathcal{O}\left(\frac{\sqrt[4]{R_\mu(f')^2 + R_\mu(f')}}{\sqrt{m\lambda^{0.75}}}\right), \quad (5)$$

which is slightly weaker than originally claimed,

$$R_\mu(f_S^{htt'}) \leq \hat{R}^{\text{loo}}(f_S^{htt'}) + \mathcal{O}\left(\sqrt{\frac{R_\mu(f')}{m\lambda}}\right). \quad (6)$$

The result is marginally worse in dependency on regularization parameter λ and performance of the source hypothesis on the target domain, $R_\mu(f')$. However, corrected bound preserves the important message for the scenario of transfer learning: for a stable algorithm, good performance of the source hypothesis on the target domain accelerates the convergence of the error measured on the training set to the expected risk.

Obtaining bound of a form (6) remains an open problem.

References

- Bousquet, O. and Elisseeff, A. Stability and Generalization. *Journal of Machine Learning Research*, 2: 499–526, 2002.
- De Vito, E., Caponnetto, A., and Rosasco, L. Model Selection for Regularized Least-Squares Algorithm in Learning Theory. *Found. Comput. Math.*, 5(1): 59–85, February 2005.
- Kuzborskij, I. and Orabona, F. Stability and hypothesis transfer learning. In *Proceedings of the International Conference on Machine Learning*, 2013.
- Petersen, K.B. and Pedersen, M.S. The matrix cookbook. *Technical University of Denmark*, 2008.