# CONSISTENT REGRESSION OF BIOPHYSICAL PARAMETERS WITH KERNEL METHODS

*Emiliano Díaz, Adrián Pérez-Suay, Valero Laparra, Gustau Camps-Valls*

Image Processing Lab (IPL), Universitat de València, València, Spain

## ABSTRACT

This paper introduces a novel statistical regression framework that allows the incorporation of consistency constraints. A linear and nonlinear (kernel-based) formulation are introduced, and both imply closed-form analytical solutions. The models exploit all the information from a set of drivers while being maximally independent of a set of auxiliary, protected variables. We successfully illustrate the performance in the estimation of chlorophyll content.

***Index Terms*—** kernel methods, regression, model inversion, consistency, vegetation monitoring

## 1. INTRODUCTION

Recent years have witnessed a successful adoption of statistical methods for model inversion, emulation and biogeophysical parameter retrieval [1]. Machine learning algorithms are flexible non-parametric models that fit the observations using large heterogeneous data. Machine learning models for parameter retrieval avoid complicated assumptions, provide fast and accurate estimates, and learn the complex relations directly from data.

Current operational vegetation products, like leaf area index (LAI), are typically produced with neural networks, Gross Primary Production (GPP) –as the largest global $CO_2$ flux driving several ecosystem functions– is estimated using ensembles of random forests, kernel methods and neural networks [2], biomass has been estimated with stepwise multiple regression [3], partial least squares regression is used for mapping canopy nitrogen [4, 5], support vector regression [6] showed high efficiency in modelling LAI, fCOVER and evapotranspiration [7, 8], and kernel methods in general [9, 10], and Gaussian Processes (GPs) in particular [11], recently provided excellent results in chlorophyll content estimation among other vegetation parameters [12, 13].

There is however an important issue that is often disregarded: statistical models learn input-output mappings from data but very often do not respect the most elemental rules of physics. They often come up with accurate, yet inconsistent predictions. For example, bio-geo-physical parameters are typically estimated with individual, indepdent models which

are typically trained separately. This common approach ignores the (potentially nonlinear) cross-relations among variables. Constraining the estimation problem is known in machine learning as *structured-output learning*, and is tightly related to *multitask learning* [14, 15]. Extension of such models to the regression setting is far from trivial, as the number of constraints increases cubically with the number of samples and outputs. Furthermore, including constraints in the regression models goes beyond consistency of model outputs; one could be interested in preserving some particular characteristics in the predictions, e.g. being independent of some ancillary information, faithful to certain variable ranges, or disregarding some information from the input (spectral) bands, just to name a few. Our notion of *consistency* is broad: we posit that a prediction is fully consistent with respect to some sensitive features if and only if the model's predictions are statistically independent of them.

In this work, we introduce a novel statistical regression framework that allows one to incorporate such broad consistency constraints. The framework builds upon [16, 17] to minimize a functional that tries to jointly minimize the empirical error and maximize the dependence of the predictions with respect to an external subset of predictors, observations or ancillary information here called *sensitive* features. Two models are derived: a linear and a nonlinear, kernel-based, consistent regression model. The new *consistency term* trades-off accuracy for consistency, and translates into an extra regularization term that can be easily interpreted. Interestingly, the proposed models come in closed-form analytical solutions, and are very easy to implement, involving only matrix inversions.

The remainder of the paper is organized as follows. Section 2 introduces notation and reviews the consistent regression models proposed. Section 3 gives experimental evidence of performance in chlorophyll content estimation. Finally, Section 4 concludes the paper with some summarizing remarks.

## 2. PROPOSED METHODOLOGY

This section starts by defining the notation and the concept of consistent regression. The proposed framework for performing consistent regression learning based on cross-covariance operators for dependence estimation in Hilbert spaces is then introduced.

## 2.1. Notation and the regularization framework

We are given $n$ samples of a response (or target) data matrix $\mathbf{Y} \in \mathbb{R}^{n \times c}$, and $d + q$ prediction variables: $d$ driver variables $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $q$ sensitive $\mathbf{S} \in \mathbb{R}^{n \times q}$. The goal is to obtain a generic prediction function (or model) $f$ for the target variable $\mathbf{Y}$ from the input data, $(\mathbf{X}, \mathbf{S})$. The goal in our framework of *consistent learning* is to predict $\mathbf{Y}$ while being maximally independent of $\mathbf{S}$.

A prediction is said to be totally consistent with respect to the sensitive features $\mathbf{S}$ if and only if $\widehat{\mathbf{Y}} \perp \mathbf{S}$. Therefore, two main ingredients are needed to perform consistent predictions: we need to ensure independence of the predictions on the sensitive variables, and simultaneously to obtain a good approximation of the target variables.

The proposed function $f$ tries to learn the relation between observed input-output data pairs $(\mathbf{x}_1, \mathbf{y}_1), \ldots, (\mathbf{x}_n, \mathbf{y}_n) \in \mathcal{X} \times \mathcal{Y}$ such that it generalizes well (good predictions $\hat{\mathbf{y}}_* = f(\mathbf{x}_*) \in \mathcal{Y}$ for the unseen input data point $\mathbf{x}_* \in \mathcal{X}$), and the predictions should be as independent as possible of the sensitive features (variables, auxiliary information or even observations). As such, the following functional should be optimized:

$$\mathcal{L} = \frac{1}{n} \sum_{i=1}^{n} V(f(\mathbf{x}_i), \mathbf{y}_i) + \lambda \, \Omega(\|f\|_{\mathcal{H}}) + \mu \, I(f(\mathbf{x}), \mathbf{s}), \quad (1)$$

where $V$ is the error cost function, $\Omega(\|f\|_{\mathcal{H}})$ acts as a regularizer of the predictive function and controls the smoothness and complexity of the model, and $I(f(\mathbf{x}), \mathbf{s})$ measures the independence between the model's predictions and the protected variables. Note that one aims to minimize the amount of information that the model shares with the sensitive variables while controlling the trade-off between fitting and independence through hyperparameters $\lambda$ and $\mu$. By setting $\mu = 0$ one obtains the ordinary (Tikhonov's regularized) functional, and by setting $\lambda = 0$ one obtains the unregularized versions of this framework.

The framework admits many variants depending on the cost function $V$, regularizer $\Omega$ and the independence measure, $I$. For example, in [18], the function $f$ was the logistic regression classifier and $I$ was a simplification of the mutual information estimate. Despite the good results reported in [18], these choices do not allow one to solve the problem in closed-form, nor to cope with more than one sensitive variable at a time, since the proposed mutual information is an uni-dimensional dependence measure. In the following section, we elaborate on this framework by using the concept of cross-covariance operators in Hilbert spaces, which lead to closed-form solutions and permit one to deal with several sensitive variables simultaneously.

## 2.2. Consistent Linear Regression

Let us now provide a straightforward instantiation of the proposed framework for consistent linear regression (CLR). We will adopt a linear predictive model for $f$, i.e. the matrix of predictions for a test data matrix $\mathbf{X}_*$ is given by $\hat{\mathbf{Y}}_* = \mathbf{X}_* \mathbf{W}$, the mean square error for the cost function $V = \|\mathbf{Y} - \mathbf{X}\mathbf{W}\|_2^2$ and the standard $\ell_2$ regularization for model weights $\Omega := \|\mathbf{W}\|_2^2$. Other choices could be made, leading to alternative formulations. In order to measure dependence, we will rely on the cross-covariance operator between the predictions and the sensitive variables in Hilbert space. Let us consider two spaces $\mathcal{Y} \subseteq \mathbb{R}^c$ and $\mathcal{S} \subseteq \mathbb{R}^q$, where random variables $(\hat{\mathbf{y}}, \mathbf{s})$ are sampled from the joint distribution $\mathbb{P}_{\mathbf{ys}}$. Given a set of pairs $\mathcal{D} = \{(\hat{\mathbf{y}}_1, \mathbf{s}_1), \ldots, (\hat{\mathbf{y}}_n, \mathbf{s}_n)\}$ of size $n$ drawn from $\mathbb{P}_{\mathbf{ys}}$, an empirical estimator of HSIC [19] allows us to define

$$I := \mathrm{HSIC}(\mathcal{Y}, \mathcal{S}, \mathbb{P}_{\mathbf{ys}}) = \|\mathbf{C}_{ys}\|_{\mathrm{HS}}^2 = \frac{1}{n^2} \mathrm{Tr}(\tilde{\mathbf{Y}}^\top \tilde{\mathbf{S}} \tilde{\mathbf{S}}^\top \tilde{\mathbf{Y}}),$$

where $\| \cdot \|_{\mathrm{HS}}$ is the Hilbert-Schmidt norm, $\mathbf{C}_{ys}$ is the empirical cross-covariance matrix between predictions and sensitive variables[1], $\tilde{\mathbf{Y}}$ and $\tilde{\mathbf{S}}$ represent the feature-centered $\mathbf{Y}$ and $\mathbf{S}$ respectively, and $\mathrm{Tr}(\cdot)$ denotes the trace operation. We want to stress that HSIC allows us to estimate dependencies between multidimensional variables, and that HSIC is zero if an only if there is no second-order dependence between $\hat{\mathbf{y}}$ and $\mathbf{s}$. In the next section we extend the formulation to higher-order dependencies with the use of kernels [9, 20].

Plugging these definitions of $f$, $V$, $\Omega$ and $I$ in Eq. (1), one can easily show that the solution has the following closed-form solution for weight estimates
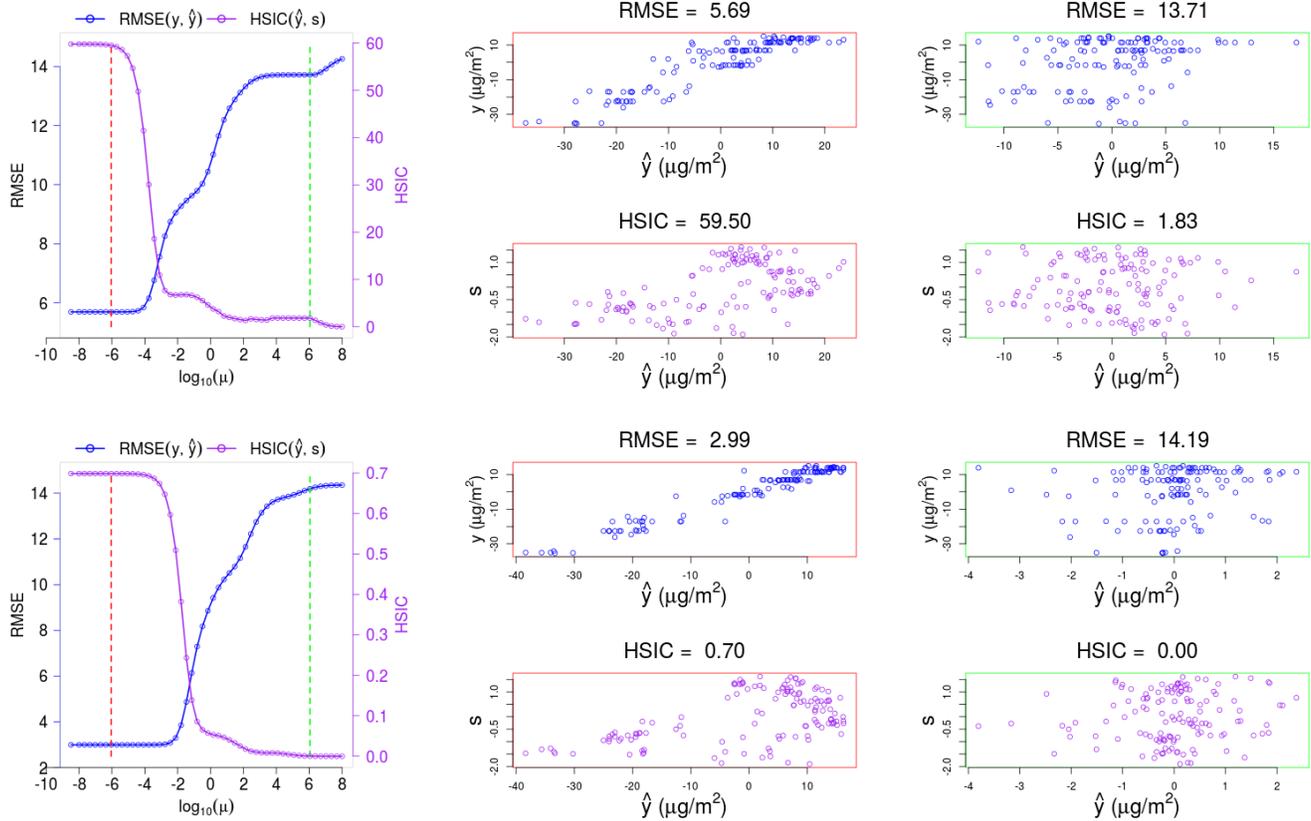
$$\widehat{\mathbf{W}} = (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} + \lambda \, \mathbf{I} + \frac{\mu}{n^2} \, \tilde{\mathbf{X}}^\top \tilde{\mathbf{S}} \tilde{\mathbf{S}}^\top \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^\top \mathbf{Y}, \quad (2)$$

where consistency is trivially controlled with $\mu$, which acts as an additional regularization term. Also note that when $\mu = 0$ the ordinary (regularized) least squares solution is obtained.

## 2.3. Consistent Kernel Regression

Let us now extend the previous model to the nonlinear case in terms of the prediction function, the regularizer and the dependence measure by means of reproducing kernels [9, 20]. We call this method the consistent kernel regression (CKR) model. We proceed in the standard way in kernel machines by mapping data $\mathbf{X}$ and $\mathbf{S}$ to a Hilbert space $\mathcal{H}$ via the mapping functions $\boldsymbol{\phi}(\cdot)$ and $\boldsymbol{\psi}(\cdot)$ respectively. This yields $\boldsymbol{\Phi}, \boldsymbol{\Psi} \in \mathcal{H} \subseteq \mathbb{R}^{d_{\mathcal{H}}}$, where $d_{\mathcal{H}}$ is the (unknown and possibly infinite) dimensionality of mapped points in $\mathcal{H}$. The corresponding kernel matrices can be defined as: $\tilde{\mathbf{K}} = \tilde{\boldsymbol{\Phi}} \tilde{\boldsymbol{\Phi}}^\top$ and $\tilde{\mathbf{K}}_S = \tilde{\boldsymbol{\Psi}} \tilde{\boldsymbol{\Psi}}^\top$. Now the prediction function is $\hat{\mathbf{Y}} = \boldsymbol{\Phi} \mathbf{W}_{\mathcal{H}}$, the regularizer is $\Omega := \|\mathbf{W}_{\mathcal{H}}\|_2^2$, and the dependence measure $I$ is the HSIC estimate between predictions $\hat{\mathbf{Y}}$ and sensitive variables $\mathbf{S}$, which can now be estimated in Hilbert spaces: $I := \mathrm{HSIC}(\mathcal{Y}, \mathcal{H}, \mathbb{P}_{\mathbf{ys}}) = \|\mathbf{C}_{ys}\|_{\mathrm{HS}}^2$. Now, by plugging all

---

[1]The covariance matrix is $\mathcal{C}_{\mathbf{ys}} = \mathbb{E}_{\mathbf{ys}}(\mathbf{ys}^\top) - \mathbb{E}_{\mathbf{y}}(\mathbf{y})\mathbb{E}_{\mathbf{s}}(\mathbf{s}^\top)$, where $\mathbb{E}_{\mathbf{ys}}$ is the expectation with respect to $\mathbb{P}_{\mathbf{ys}}$, and $\mathbb{E}_{\mathbf{y}}$ is the marginal expectation with respect to $\mathbb{P}_{\mathbf{y}}$ (hereafter we assume that all these quantities exist).

**Fig. 1**. Evolution of the RMSE [$\mu$g/m$^2$] and HSIC for predicting Chl-a with either linear (top) or kernel (bottom) regression as a function of the consistency parameter $\mu$. Scatter plots illustrate the accuracy (blue points, predicted vs. observed chlorophyll content) and consistency (purple points, predicted chlorophyll content vs. sensitive band) for two choices of $\mu$ (low and high consistency correspond to red and green respectively).

these terms in the cost function, using the representer's theorem $\mathbf{W}_{\mathcal{H}} = \tilde{\boldsymbol{\Phi}}^{\top}\boldsymbol{\Lambda}$ and after some simple linear algebra, we obtain the dual weights in closed-form

$$\boldsymbol{\Lambda} = (\tilde{\mathbf{K}} + \lambda\mathbf{I} + \frac{\mu}{n^2}\tilde{\mathbf{K}}_S\tilde{\mathbf{K}})^{-1}\mathbf{Y}, \qquad (3)$$

which can be used for prediction with a new point $\mathbf{x}_*$ by using $\hat{\mathbf{y}}_* = \mathbf{k}_*\boldsymbol{\Lambda}$, where $\mathbf{k}_* = [K(\mathbf{x}_*, \mathbf{x}_1), \ldots, K(\mathbf{x}_*, \mathbf{x}_n)]^{\top}$. Note that in the case where $\mu = 0$ the method reduces to standard kernel ridge regression (KRR) method [9]. Note that centering points in feature spaces can be done implicitly with kernels [9]: a kernel matrix $\mathbf{K}$ is centered by doing $\tilde{\mathbf{K}} = \mathbf{H}\mathbf{K}\mathbf{H}$, where $\mathbf{H} = \mathbf{I} - \frac{1}{n}\mathbb{1}\mathbb{1}^{\top}$.

## 3. EXPERIMENTS

This section presents the results of the application of our consistent regression framework to a remote sensing problem. In particular, we illustrate the performance of the proposed method to retrieve consistent estimates of chlorophyll content from hyperspectral images.

### 3.1. Data collection

The data were obtained in the SPARC-2003 (SPectra bARrax Campaign) and SPARC-2004 campaigns in Barrax, Spain. The region consists of approximately 65% dry land and 35% irrigated land. The methodology applied to obtain the *in situ* leaf-level Chl$_{ab}$ data consisted of measuring samples with a calibrated CCM-200 Chlorophyll Content Meter in the field. Chl measurements were between 2 and 55 $\mu$g/cm$^2$. Additionally, 30 random bare soil spectra with zero chlorophyll value were added to broaden the dataset to non-vegetated samples. Concurrently, we used CHRIS images Mode 1 (62 spectral bands, 34m spatial resolution at nadir). The images were pre-processed, geometrically and atmospherically corrected. A total of $n = 136$ datapoints in a 62-dimensional space and the measured chlorophyll concentration constitute the database.

### 3.2. Results

The experiment deals with the prediction of chlorophyll content while forcing the model to be as independent as possible

from the bands beyond the NIR. It is physically understood that the chlorophyll content drives reflectance mostly in the red edge. Hence, our sensitive variables are the channels far beyond the NIR spectrum. Figure 1 shows the results that were obtained. We show the root mean square error (RMSE) and the HSIC for different levels of the consistency parameter $\mu$ (and, implicitly, the corresponding optimal $\lambda(\mu)$) for both the linear and non-linear models.

Several conclusions can be obtained. First, one can readily recognize, in both cases, a trade-off between obtaining accurate predictions and imposing consistency. As we increase the consistency parameter $\mu$ we obtain predictions with greater independence to the sensitive variables but the accuracy of the prediction deteriorates. Second, although more accurate predictions can be obtained with the non-linear, kernel model (RMSE=2.99 mg/m$^3$ versus RMSE=5.69 mg/m$^3$), the rate of deterioration of the accuracy when consistency is imposed is greater for the non-linear, kernel model than for its linear counterpart. Figure 1 also illustrates the quality of predictions and consistency for two choices of the consistency parameter $\mu$. It can be noted that good models in terms of accuracy (red point) lead to more correlated predictions with the sensitive bands, while enforcing the constraints (blue point) lead to higher independence but poor fitting results for both linear (RMSE=13.71 mg/m$^3$) or nonlinear (RMSE=14.19 mg/m$^3$) models.

## 4. CONCLUSIONS

This paper presented a novel statistical regression framework that allows one to incorporate consistency constraints. The methodology confers, to both linear and nonlinear statistical regression, methods whose solution can be expressed in closed-form. The models exploit all the information from a set of covariates while being maximally independent of a set of auxiliary, protected variables. We successfully illustrated the performance for the estimation of chlorophyll content while being independent to particular spectral bands.

## 5. REFERENCES

[1] G. Camps-Valls, D. Tuia, L. Gómez-Chova, and J. Malo, Eds., *Remote Sensing Image Processing*, Morgan & Claypool, Sept 2011.

[2] G. Tramontana, M. Jung, G. Camps-Valls, K. Ichii, B. Raduly, M. Reichstein, C. R. Schwalm, M. A. Arain, A. Cescatti, G. Kiely, L. Merbold, P. Serrano-Ortiz, S. Sickert, S. Wolf, and D. Papale, "Predicting carbon dioxide and energy fluxes across global fluxnet sites with regression algorithms," *Biogeosciences Discussions*, vol. 2016, pp. 1–33, 2016.

[3] L. R. Sarker and J. E. Nichol, "Improved forest biomass estimates using ALOS AVNIR-2 texture indices," *Rem. Sens. Env.*, vol. 115, no. 4, pp. 968–977, 2011.

[4] N. C. Coops, M-L. Smith, M.E. Martin, and S. V. Ollinger, "Prediction of eucalypt foliage nitrogen content from satellite-derived hyperspectral data," *IEEE Trans. Geosc. Rem. Sens.*, vol. 41, no. 6, pp. 1338–1346, Jun 2003.

[5] P.A. Townsend, J.R. Foster, R.A. Jr. Chastain, and W.S. Currie, "Application of imaging spectroscopy to mapping canopy nitrogen in the forests of the central Appalachian Mountains using Hyperion and AVIRIS," *IEEE Trans. Geosc. Rem. Sens.*, vol. 41, no. 6, pp. 1347–1354, June 2003.

[6] A. J. Smola and B. Schölkopf, "A tutorial on support vector regression," *Statistics and Computing*, vol. 14, pp. 199–222, 2004.

[7] S.S. Durbha, R.L. King, and N.H. Younan, "Support vector machines regression for retrieval of leaf area index from multiangle imaging spectroradiometer," *Rem. Sens. Env.*, vol. 107, no. 1-2, pp. 348–361, 2007.

[8] F. Yang, M.A. White, A.R. Michaelis, K. Ichii, H. Hashimoto, P. Votava, A-Xing Zhu, and R.R. Nemani, "Prediction of Continental-Scale Evapotranspiration by Combining MODIS and AmeriFlux Data Through Support Vector Machine," *IEEE Trans. Geosc. Rem. Sens.*, vol. 44, no. 11, pp. 3452–3461, nov. 2006.

[9] J. Shawe-Taylor and N. Cristianini, *Kernel Methods for Pattern Analysis*, Cambridge University Press, 2004.

[10] G. Camps-Valls and L. Bruzzone, Eds., *Kernel methods for Remote Sensing Data Analysis*, Wiley & Sons, UK, Dec 2009.

[11] C. E. Rasmussen and C. K. I. Williams, *Gaussian Processes for Machine Learning*, The MIT Press, New York, 2006.

[12] R. Furfaro, R. D. Morris, A. Kottas, M. Taddy, and B. D. Ganapol, "A Gaussian Process Approach to Quantifying the Uncertainty of Vegetation Parameters from Remote Sensing Observations," *AGU Fall Meeting Abstracts*, pp. A261+, Dec 2006.

[13] G. Camps-Valls, J. Verrelst, J. Muñoz-Marí, V. Laparra, F. Mateo-Jimenez, and J. Gomez-Dans, "A survey on gaussian processes for earth observation data analysis: A comprehensive investigation," *IEEE Geoscience and Remote Sensing Magazine*, , no. 6, June 2016.

[14] J. Leiva, L. Gómez-Chova, and G. Camps-Valls, "Multitask remote sensing data classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 50, Oct 2012.

[15] Devis Tuia, Jordi Muñoz-Marí, Mikhail F. Kanevski, and Gustavo Camps-Valls, "Structured output SVM for remote sensing image classification," *Signal Processing Systems*, vol. 65, no. 3, pp. 301–310, 2011.

[16] Toshihiro Kamishima, Shotaro Akaho, Hideki Asoh, and Jun Sakuma, "The independence of fairness-aware classifiers," *2013 IEEE 13th International Conference on Data Mining Workshops*, vol. 00, pp. 849–858, 2013.

[17] A. Pérez-Suay, V. Laparra, G. Mateo-García, J. Mu ñ oz Marí, L. Gómez-Chova, and G. Camps-Valls, "Fair kernel learning," in *European Conference on Machine Learning (ECML)*, Skopje, Macedonia, 18-22 September 2017.

[18] Toshihiro Kamishima, Shotaro Akaho, Hideki Asoh, and Jun Sakuma, *Fairness-Aware Classifier with Prejudice Remover Regularizer*, pp. 35–50, Springer Berlin Heidelberg, Berlin, Heidelberg, 2012.

[19] A. Gretton, R. Herbrich, and A. Hyvärinen, "Kernel methods for measuring independence," *Journal of Machine Learning Research*, vol. 6, pp. 2075–2129, 2005.

[20] B. Schölkopf and A. Smola, *Learning with Kernels – Support Vector Machines, Regularization, Optimization and Beyond*, MIT Press Series, 2002.