

A New Look at an Old Problem: A Universal Learning Approach to Linear Regression

Koby Bibas
School of Electrical Engineering
Tel Aviv University
Email: kobybibas@gmail.com

Yaniv Fogel
School of Electrical Engineering
Tel Aviv University
Email: Yaniv.fogel8@gmail.com

Meir Feder
School of Electrical Engineering
Tel Aviv University
Email: meir@eng.tau.ac.il

Abstract—Linear regression is a classical paradigm in statistics. A new look at it is provided via the lens of universal learning. In applying universal learning to linear regression the hypotheses class represents the label $y \in \mathcal{R}$ as a linear combination of the feature vector $x^T \theta$ where $x \in \mathcal{R}^M$, within a Gaussian error. The Predictive Normalized Maximum Likelihood (pNML) solution for universal learning of individual data can be expressed analytically in this case, as well as its associated learnability measure. Interestingly, the situation where the number of parameters M may even be larger than the number of training samples N can be examined. As expected, in this case learnability cannot be attained in every situation; nevertheless, if the test vector resides mostly in a subspace spanned by the eigenvectors associated with the large eigenvalues of the empirical correlation matrix of the training data, linear regression can generalize despite the fact that it uses an “over-parametrized” model. We demonstrate the results with a simulation of fitting a polynomial to data with a possibly large polynomial degree.

I. INTRODUCTION

Linear regression, using least squares, is probably one of the most standard techniques in statistics, [1]. This work provides a new view for this problem based on recent results in universal learning. In particular, the common assumption in linear regression is that the number of training samples needs to be much higher than the number of features in order to be able to generalize [2]. Recently, the success of Deep Neural Network (DNN) in which the number of learnable parameters is in several order of magnitudes greater than the size of the feature space requires rethinking that assumption. The new view we provide will show that sometimes generalization can be attained even in the “over-parameterized” regime.

Before diving into this analysis, a short introduction to universal learning is provided. In the common situation of supervised machine learning, a training set is given consisting of N pairs $z^N = \{(x_i, y_i)\}_{i=1}^N$, where $x \in \mathcal{X}$ is the data or the features and $y \in \mathcal{Y}$ is the label. Then, a new x is given and the task is to predict its corresponding label y . In the information theoretic framework considered in a variety of works, e.g., [3] and more recently [4], prediction is done by assigning a probability distribution $q(\cdot|x)$ to the unknown label, and the prediction loss is the log-loss:

$$\mathcal{L}(q; x, y) = -\log q(y|x). \quad (1)$$

Clearly a reasonable goal is to find the predictor q with the minimal loss for the test sample. However, this problem is ill-posed unless additional assumptions are made.

First, a “model” class, or “hypotheses” class must be defined. This class is a set of conditional probability distributions

$$P_\Theta = \{p_\theta(y|x), \theta \in \Theta\} \quad (2)$$

where Θ is a general index set. This is equivalent to saying that there is a set of stochastic functions $\{y = g_\theta(x), \theta \in \Theta\}$ used to explain the relation between x and y .

Next, assumptions must be made on how the data and the labels are generated. In the stochastic setting it is assumed that there is a true probabilistic relation between x and y given by an (unknown) model from the class P_Θ . A more general setting, used in the variety of works in machine learning, is the Probably Approximately Correct (PAC) established in [5]. In PAC x and y are assumed to be generated by some source $P(x, y) = P(x)P(y|x)$, but unlike the standard stochastic setting $P(y|x)$ is not necessarily a member of the hypothesis class. In both the stochastic and PAC settings the goal is to perform as well as a learner that knows the true probability.

The most general setting, however, and the one used in this paper is the individual, where the data and labels of both the training and test are specific and individual. In this setting the goal can no longer be to perform as well as a learner that knows the true probability. Instead, following [3], the goal is to seek a learner that can compete with a “genie” or a reference learner that knows the desired label value, but is restricted to use a model from the given hypotheses class P_Θ . In addition, as discussed in [4], the reference has no knowledge which of the samples is the test. Thus, the reference chooses

$$\hat{\theta}(z^N, x, y) = \arg \max_{\theta} [p_\theta(y|x) \cdot \prod_{i=1}^N p_\theta(y_i|x_i)] \quad (3)$$

The log-loss difference between a universal learner q and the reference is the regret:

$$R(z^N, x, y, q) = \log \frac{p_{\hat{\theta}(z^N, x, y)}(y|x)}{q(y|x; z^N)}. \quad (4)$$

As advocated in [4], the chosen universal learner solves:

$$\min_q \max_y R(z^N, x, y, q) = R^*(z^N, x) \quad (5)$$

Following [6] this learner, called the Predictive Normalized Maximum Likelihood (pNML), is given by

$$q_{\text{pNML}}(y|x; z^N) = \frac{p_{\hat{\theta}(z^N, x, y)}(y|x)}{\sum_{y \in \mathcal{Y}} p_{\hat{\theta}(z^N, x, y)}(y|x)}. \quad (6)$$

Note that this pNML probability assignment was essentially proposed earlier, see [7], [8], with a different motivation as one of the possible variations of the Normalized Maximum Likelihood (NML) method of [6] for universal prediction.

In order to obtain the pNML the following procedure is executed: assuming the label of the test data is known. Find the best model that fits it with the training samples, and predict the assumed label by it. Repeat the process for all possible labels. Then, normalize to get a valid probability distribution which is the pNML learner. The regret of the pNML, $R^*(z^N, x)$ is the logarithm of its normalization factor

$$R^*(z^N, x) = \log \left\{ \sum_{y \in \mathcal{Y}} p_{\hat{\theta}(z^N, x, y)}(y|x) \right\}. \quad (7)$$

In considering linear regression, $y \in \mathcal{R}$ is the scalar label, $x \in \mathcal{R}^M$ is the feature vector (sometimes the first component of x is set to 1 to formulate affine linear relation), and the model class is the set:

$$\left\{ p_{\theta}(y|x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (y - x^T \theta)^2 \right\}, \theta \in \mathcal{R}^M \right\} \quad (8)$$

That is, the label y is a linear combination of the components of x , within a Gaussian noise. As shown below, in this case the pNML and its learnability measure can be evaluated explicitly.

The pNML approach deviates from the standard Empirical Risk Minimization (ERM) [9] approach. In ERM, given a training set and hypothesis class $\{p_{\theta}(y|x), \theta \in \Theta\}$, a learner that minimizes the loss over the training set is chosen:

$$q_{\text{ERM}}(y|x) = \underset{p_{\theta}}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^N \mathcal{L}(p_{\theta}; x_i, y_i). \quad (9)$$

In the linear regression model (8), one chooses the least squares solution over the training set for the linear coefficients. This, however, may lead to large generalization log-loss error.

The paper has two main contributions. First, it provides an explicit analytical solution for the pNML learner and its “learnability” measure (which is the minmax regret (7)) for the linear regression hypothesis class. This includes also the regularized case where the norm of the coefficients vector is constrained. Second, based on the analysis of the learnability measure, it is shown that even in the over-parameterized case where the number of parameters M may be larger than the training size N , if the test data comes from a “learnable space” successful generalization occurs. This phenomenon may explain why other over-parameterized models such as deep neural networks are successful for “learnable” data.

The paper outline is as follows. Section II presents some related work. Section III provides the formal problem definition, while the pNML evaluation for the regression problem is

given in sections IV and V. In depth analysis of the learnable space is given in section VI. Simulation of the pNML and its regret for the problem of fitting a polynomial to the data is described in section VII and conclusion in section VIII.

II. RELATED WORKS

In this section we briefly mention related work on model generalization and least squares regression.

Model Generalization. Understating the model generalization capabilities is considered a fundamental problem in machine learning [10]. As noted, most of the theoretical work in learning use the PAC setting. In that setting, a common measure is the VC Dimension that can be used to upper bound on the test generalization error. For DNN’s, the VC dimension is linear with the number of parameters [11], yet the empirical evidence demonstrates that DNN’s have state of the art generalization performance. This makes the VC dimension irrelevant for assessing the generalization error of DNN’s.

Least Squares The least squares algorithm is widely used in linear regression due to its robust performance and simplicity of implementation. In addition to the explicit formula for its solution, it can be solved sequentially, via the Recursive Least Squares (RLS) algorithm, which is an efficient online method for finding the linear predictor that minimizes the squared error over the training data [12]. In learning linear regression models, the least squares is used for prediction in the ERM approach and in a Bayesian linear regression approach [13]. This paper provides an analysis in the individual setting, in the realm of the pNML approach.

III. LINEAR REGRESSION: FORMAL PROBLEM DEFINITION

Given N pairs of data and labels $\{x_i, y_i\}_{i=1}^N$ where $x_i \in \mathcal{R}^M$, $y_i \in \mathcal{R}$ are deterministic, the model takes the form:

$$\begin{aligned} y_1 &= x_1^T \theta + e_1 \\ &\vdots \\ y_N &= x_N^T \theta + e_N \end{aligned} \quad (10)$$

where $\theta \in \mathcal{R}^M$ are the learnable parameters and the $e_i \in \mathcal{R}$ are zero mean, Gaussian, independent with variance of σ^2 . The goal is to predict y based on a new data sample x . Under the assumptions y , conditioned on x , has a normal distribution that depends on the learnable parameters θ

$$p_{\theta}(y) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (y - x^T \theta)^2 \right\} \quad (11)$$

The unknown parameter vector θ belongs to a set Θ , which in the general case is the entire \mathcal{R}^M . In the regularized version (that leads to Ridge regression [14]), Θ is the sphere $|\theta| \leq A$. In the next section the pNML will be evaluated for this hypotheses class. Recall that the pNML learner of y given the test sample x and the training set $z^N = \{(x_i, y_i)\}_{i=1}^N$ is given by:

$$q_{\text{pNML}}(y|x; z^N) = \frac{1}{\Gamma} p_{\hat{\theta}(z^N, x, y)}(y|x). \quad (12)$$

where in the linear regression case

$$\hat{\theta}(z^N, x, y) = \arg \min_{\theta \in \Theta} \left[\sum_{i=1}^N (y_i - x_i^T \theta)^2 + (y - x^T \theta)^2 \right] \quad (13)$$

and where Γ is the the normalization factor:

$$\Gamma(z^N, x) = \int_R p_{\hat{\theta}(z^N, x, y)}(y|x) dy, \quad (14)$$

The goal is to find an analytic expression for (12) and for the learnability measure $\log \Gamma(z^N, x)$, the minmax regret value.

IV. PNML EVALUATION

The following notation is used. $X \in R^{M \times N+1}$ is the matrix which contains all the training data along with the test sample and $Y \in R^{N+1}$ is the vector which contains all the labels including the test label, i.e.,

$$X = \begin{bmatrix} x_1 & \dots & x_N & x \end{bmatrix}, \quad Y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \\ y \end{bmatrix} \quad (15)$$

Assuming that the test label y is given, the optimal solution under least squares:

$$\hat{\theta}(z^N, x, y) = \theta_{N+1}^* = (X^T X)^{-1} X^T Y \quad (16)$$

By the recursive least square (RLS) formulation [12]:

$$\theta_{N+1}^* = \theta_N^* + P_{N+1} x (y - \hat{y}) \quad (17)$$

where $\hat{y} = x^T \theta_N^*$ is the ERM prediction based on the samples $\{(x_i, y_i)\}_{i=1}^N$ and

$$P_{N+1} = (X X^T)^{-1}. \quad (18)$$

Note that in RLS P_{N+1} is also calculated recursively from P_N , but this is not needed at this point. Now,

$$\begin{aligned} P_{\theta_{N+1}^*}(y) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (y - x^T \theta_{N+1}^*)^2 \right\} = \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (y - x^T (\theta_N^* + \right. \\ &\quad \left. P_{N+1} x (y - \hat{y}))^2 \right\} = \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(1 - x^T P_{N+1} x)^2}{2\sigma^2} (y - \hat{y})^2 \right\}. \end{aligned} \quad (19)$$

To get the pNML normalization factor (14), we integrate

$$\begin{aligned} \Gamma &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(1 - x^T P_{N+1} x)^2}{2\sigma^2} (y - \hat{y})^2 \right\} dy \\ &= \frac{1}{1 - x^T P_{N+1} x} = \frac{1}{1 - x^T (X X^T)^{-1} x} \end{aligned} \quad (20)$$

Thus, the pNML distribution of y given x is:

$$q_{\text{pNML}}(y|x; z^N) = \frac{1}{\Gamma} p_{\theta_{N+1}^*}(y|x) = \frac{1 - x^T P_{N+1} x}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(1 - x^T P_{N+1} x)^2}{2\sigma^2} (y - \hat{y})^2 \right\} \quad (21)$$

and its associate learnability measure or regret:

$$\log \Gamma(z^N, x) = \log \left(\frac{1}{1 - x^T (X X^T)^{-1} x} \right). \quad (22)$$

V. PNML WITH REGULARIZATION

Assuming now that the model class Θ is constrained to the sphere $|\theta| \leq A$, for some A . Using a Lagrange multiplier λ we get the Tikhonov regularization (or Ridge regression), where the expression to minimize is now:

$$\mathcal{L}(z^N) = \sum_{i=1}^N |y_i - x_i^T \theta|^2 + \lambda |\theta|^2 \quad (23)$$

With the test data, the “regularized” least square solution is:

$$\hat{\theta}(z^N, x, y) = \theta_{N+1}^* = (X^T X + \lambda I)^{-1} X^T Y \quad (24)$$

Here too the RLS formula holds:

$$\theta_{N+1}^* = \theta_N^* + P_{N+1} x (y - \hat{y}) \quad (25)$$

However, now

$$P_{N+1} = (X^T X + \lambda I)^{-1}. \quad (26)$$

The rest of the evaluation can follow as before. Thus, the pNML learner in this “regularized” case is

$$\begin{aligned} q_{\text{pNML}}(y|x; z^N, \lambda) &= \frac{1 - x^T (X X^T + \lambda I)^{-1} x}{\sqrt{2\pi\sigma^2}} \\ &\cdot \exp \left\{ -\frac{(1 - x^T (X X^T + \lambda I)^{-1} x)^2}{2\sigma^2} (y - \hat{y})^2 \right\} \end{aligned} \quad (27)$$

and the associated regret or the log-normalization factor:

$$\log \Gamma(z^N, x) = \log \left(\frac{1}{1 - x^T (X X^T + \lambda I)^{-1} x} \right) \quad (28)$$

Note that regularization can help in the case where $X X^T$, the un-normalized correlation matrix of the data is ill conditioned. In the next section we find the “learnable space” for the linear regression problem and observe situations where this regularization is needed.

VI. LEARNABLE SPACE

In order to understand for which test sample the trained model generalizes well we need to look at the regret expression (22). High regret means that the pNML learner is very far from the genie and therefore we may not trust its predictions. Low regret, on the other hand, means the model is as good as a genie who knows the true test label, and so it is trusted.

Consider the matrix $X_N = [x_1, x_2, \dots, x_N]$, composed of the training data, and apply the singular value decomposition (SVD) on it, i.e., $X_N = U\Sigma V^T$ with $U \in R^{M \times M}$, Σ is a rectangular diagonal matrix of the singular values and $V \in R^{N \times N}$. The expression $x^T(XX^T)^{-1}x$ can be rewritten as:

$$x^T(XX^T)^{-1}x = x^T \left([U\Sigma V^T \quad x] \begin{bmatrix} V\Sigma^T U^T \\ x^T \end{bmatrix} \right)^{-1} x = x^T (U\Sigma\Sigma^T U^T + xx^T)^{-1} x. \quad (29)$$

Denote by R_N the empirical correlation matrix of the training:

$$R_N = \frac{1}{N} U\Sigma\Sigma^T U^T = UHU^T \quad R_N^{-1} = UH^{-1}U^T \quad (30)$$

where H is a diagonal matrix with $H_{ii} = \eta_i$, the eigenvalues of R_N . By the matrix inversion lemma, see [15], we have:

$$x^T(XX^T)^{-1}x = x^T \left[\frac{1}{N} R_N^{-1} - \frac{\frac{1}{N^2} R_N^{-1} x x^T R_N^{-1}}{1 + \frac{1}{N} x^T R_N^{-1} x} \right] x. \quad (31)$$

Denote $\gamma = x^T R_N^{-1} x$. We can simplify the expression:

$$x^T(XX^T)^{-1}x = \frac{1}{N} \gamma - \frac{\frac{1}{N^2} \gamma^2}{1 + \frac{1}{N} \gamma} = \frac{\frac{1}{N} \gamma}{1 + \frac{1}{N} \gamma}. \quad (32)$$

Plugging in the regret formula (22):

$$\log \Gamma = \log \left(\frac{1}{1 - \frac{\frac{1}{N} \gamma}{1 + \frac{1}{N} \gamma}} \right) = \log \left(1 + \frac{1}{N} \gamma \right). \quad (33)$$

Let u_i be the eigenvectors of the empirical correlation matrix of the training data. Express γ by $x^T u_i$, the projections of x on u_i :

$$\gamma = [x^T u_1 \quad \dots \quad x^T u_M] \begin{bmatrix} \frac{1}{\eta_1} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \frac{1}{\eta_M} \end{bmatrix} \begin{bmatrix} u_1^T x \\ \vdots \\ u_M^T x \end{bmatrix} = \sum_{i=1}^M \frac{(x^T u_i)^2}{\eta_i}. \quad (34)$$

The final regret expression is thus:

$$\log \Gamma = \log \left(1 + \frac{1}{N} \sum_{i=0}^M \frac{(x^T u_i)^2}{\eta_i} \right). \quad (35)$$

If the test sample x lies mostly in the subspace spanned by the eigenvectors with large eigenvalues, then the model can generalize well even if $M > N$.

VII. SIMULATION

In this section we present some simulations that demonstrate the results above. We chose the problem of fitting polynomial to data, which is a special case of linear regression. The simulation show prediction and generalization capabilities in a variety of polynomial degrees and regularization factors.

In the first experiment we generated 3 random points, t_0, t_1, t_2 , uniformly in the interval $[-1, 1]$. These points are the training set and are shown in Figure 1 (top) as red dots.

The relation between y and t is given by as polynomial of degree 2. Thus, the X matrix of section III is given by:

$$X = \begin{bmatrix} 1 & 1 & 1 \\ t_0 & t_1 & t_2 \\ t_0^2 & t_1^2 & t_2^2 \end{bmatrix}. \quad (36)$$

Based on the training we predict a probability for all t values in the interval $[-1, 1]$ using (27) with a regularization factor λ of 0, 0.1 and 1.0. It is shown in Figure 1 (top) that without regularization ($\lambda = 0$), the blue curve fits the data exactly. As λ increases the fitted curve becomes less steep but tends to fit less the training data.

Figure 1 (bottom) shows the regret, given by (22), for the degree 2 polynomial model for all $t \in [-1, 1]$ where the training t_i 's are marked in red on the x axis. We can see that around the training data the regret is very low in comparison to areas where training data does not exist. In addition, with regularization, the overall regret is lower and is greater than 1.0 only on the edges of the interval. For all regularization terms, the regret values increases as moving away from the training data.

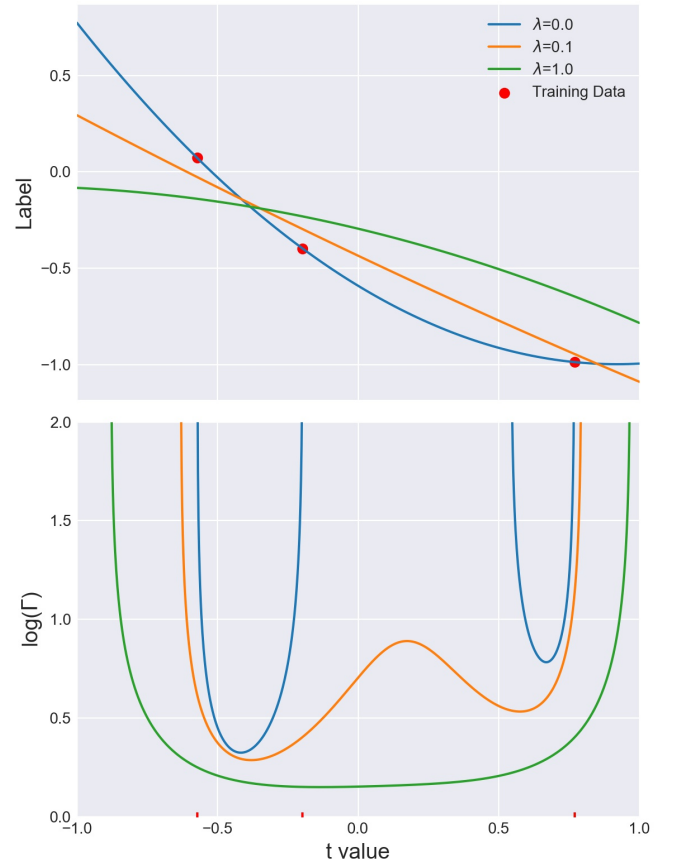


Fig. 1. **Least squares predictor with variety of regularization terms.** (Top) The least squares estimator fitted to the training data (in red) with different values of regularization term. (Bottom) The regret of the pNML learner from (22) on the interval $[-1, 1]$. The training data t values are marked in red on the x axis.

Next, we simulate the case of fitting polynomials with different degrees. Again, we generated 3 random points uniformly in the interval $[-1, 1]$. The matrix X is now:

$$X = \begin{bmatrix} 1 & 1 & 1 \\ t_0 & t_1 & t_2 \\ \vdots & \vdots & \vdots \\ t_0^{\text{Poly Deg}} & t_1^{\text{Poly Deg}} & t_2^{\text{Poly Deg}} \end{bmatrix}. \quad (37)$$

Figure 2 (top) shows the predicted label for every t value in $[-1, 1]$ for the different polynomial degrees. To avoid singularities we used the regularized version with $\lambda = 10^{-4}$. The training set is shown by red dots in the figure. Note that for a polynomial of degree 3, the number of parameters is greater than the size of the training set. Nevertheless, the prediction accuracy near the training samples is similar to that of a degree 2 polynomial. Figure 2 (bottom) shows the regret (or learnability) of the three pNML learners corresponding to model classes of polynomials with the various degrees. All the learners have regret values that are small near the training samples and large as t drifts away from these samples.

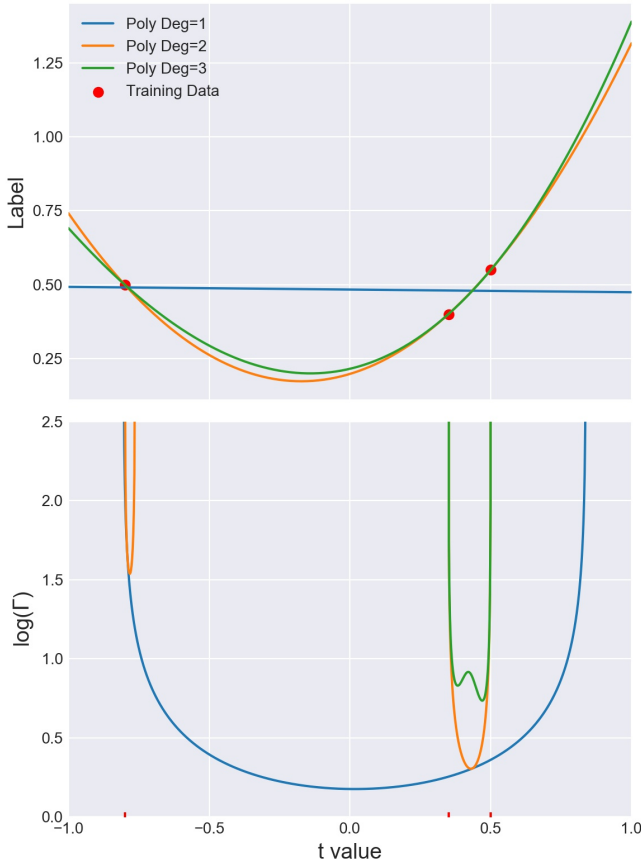


Fig. 2. **Least squares predictor with different polynomial degree.** (Top) pNML least squares predictions with different polynomial degrees. (Bottom) The regret of the pNML learners from (22) on the interval $[-1, 1]$. The training data t values are marked in red on the x axis.

VIII. CONCLUSIONS

In this paper, we provided an explicit analytical solution of the pNML universal learning scheme and its learnability measure for the linear regression hypothesis class. Interestingly, the predicted universal pNML assignment is Gaussian with a mean that is equal to that of the ERM, but with a variance that increases by a factor Γ whose logarithm is the learnability measure. Analyzing Γ we can observe the “learnability space” for this problem. Specifically, if a test sample mostly lies in the subspace spanned by the eigenvectors associated with large eigenvalues of the empirical correlation matrix then the learner can generalize well, even in an over-parameterized case where the regression dimension is larger than the number of training samples. Finally, we have provided a simulation of the pNML least squares prediction for polynomial interpolation with different regularization factors and polynomial degrees.

This work suggests a number of potential directions for future work, some are already explored in an accompanying paper [16]. We conjecture that as in linear regression other “over-parameterized” model classes are learnable at least locally at some learnable space, that can be inferred from the pNML solution. This notion is indeed corroborated by the findings in [16].

REFERENCES

- [1] C. L. Lawson and R. J. Hanson, *Solving least squares problems*. Siam, 1995, vol. 15.
- [2] G. James, D. Witten, T. Hastie, and R. Tibshirani, *An introduction to statistical learning*. Springer, 2013, vol. 112.
- [3] N. Merhav and M. Feder, “Universal prediction,” *IEEE Transactions on Information Theory*, vol. 44, no. 6, pp. 2124–2147, 1998.
- [4] Y. Fogel and M. Feder, “Universal supervised learning for individual data,” *arXiv preprint arXiv:1812.09520*, 2018.
- [5] L. G. Valiant, “A theory of the learnable,” *Communications of the ACM*, vol. 27, no. 11, pp. 1134–1142, 1984.
- [6] Y. M. Shtarkov, “Universal sequential coding of single messages,” *Problemy Peredachi Informatsii*, vol. 23, no. 3, pp. 3–17, 1987.
- [7] T. Roos and J. Rissanen, “On sequentially normalized maximum likelihood models,” 2008.
- [8] T. Roos, T. Silander, P. Kontkanen, and P. Myllymaki, “Bayesian network structure learning using factorized nml universal models,” in *Information Theory and Applications Workshop, 2008*. IEEE, 2008, pp. 272–276.
- [9] V. Vapnik, “Principles of risk minimization for learning theory,” in *Advances in neural information processing systems*, 1992, pp. 831–838.
- [10] —, *The nature of statistical learning theory*. Springer science & business media, 2013.
- [11] E. D. Sontag, “Vc dimension of neural networks,” *NATO ASI Series F Computer and Systems Sciences*, vol. 168, pp. 69–96, 1998.
- [12] M. H. Hayes, “9.4: Recursive least squares,” *Statistical Digital Signal Processing and Modeling*, p. 541, 1996.
- [13] K. W. Fornalski, “Applications of the robust bayesian regression analysis,” *International Journal of Society Systems Science*, vol. 7, no. 4, pp. 314–333, 2015.
- [14] A. Hoerl and R. Kennard, “Ridge regression: Biased estimation for nonorthogonal problems,” *Technometrics*, vol. 12, 1970.
- [15] W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery, “Section 2.7. 1 sherman–morrison formula,” *Numerical Recipes: The Art of Scientific Computing (3rd ed.)*. Cambridge University Press, New York, vol. 1, pp. 55–67, 2007.
- [16] K. Bibas, Y. Fogel, and M. Feder, “Deep pnml: Predictive noramalized maximum likelihood for deep neural networks,” *Submitted, International Conference on Machine Learning (ICML)*, 2019.