# Nonlinear Function Estimation
# with Empirical Bayes
# and Approximate Message Passing

Hangjin Liu,[†] You (Joe) Zhou,[†] Ahmad Beirami,[‡] and Dror Baron[†]

[†]Department of Electrical and Computer Engineering, North Carolina State University, Raleigh, NC 27695, USA

[‡]Research Laboratory of Electronics, Massachusetts Institute of Technology, Cambridge, MA 02139, USA

Email: {hliu25,yzhou26,barondror}@ncsu.edu, beirami@mit.edu

*Abstract*—**Nonlinear function estimation is core to modern machine learning applications. In this paper, to perform nonlinear function estimation, we reduce a nonlinear inverse problem to a linear one using a polynomial kernel expansion. These kernels increase the feature set, and may result in poorly conditioned matrices. Nonetheless, we show several examples where the matrix in our linear inverse problem contains only mild linear correlations among columns. The coefficients vector is modeled within a Bayesian setting for which approximate message passing (AMP), an algorithmic framework for signal reconstruction, offers Bayes-optimal signal reconstruction quality. While the Bayesian setting limits the scope of our work, it is a first step toward estimation of real world nonlinear functions. The coefficients vector is estimated using two AMP-based approaches, a Bayesian one and empirical Bayes. Numerical results confirm that our AMP-based approaches learn the function better than LASSO, offering markedly lower error in predicting test data.**

*Index Terms*—**Approximate message passing, function estimation, kernel regression, nonlinear functions, Taylor series.**

## I. INTRODUCTION

A pervasive trend in modern society is that ever-larger amounts of data are being collected and analyzed in order to explain various phenomena. In supervised learning, many variables (also referred to as features) that may relate to and thus help explain the phenomena of interest are observed, and the goal is to learn a function — often a nonlinear one — that relates the explanatory variables to the phenomena of interest. More specifically, we have a multivariate nonlinear function, $\mathbf{f}(\cdot)$, and collect noisy samples of it; our goal is to estimate $\mathbf{f}(\cdot)$. At its core, this is *multivariate nonlinear function estimation*; it could also be interpreted as nonlinear regression or feature selection. Algorithms for solving such problems must be robust to noisy observations and outliers, backed up by fundamental mathematical analysis, support missing data, and have a fast implementation that scales well to large-scale problems. Such algorithms will impact many disciplines, such as health informatics [1], social networks, and finance [2–4].

**Example applications:** Let us describe how nonlinear function estimation can be used in financial prediction [2–4]. A typical approach to estimate expected returns uses a linear factor model, which is tuned to work well on training data, $y_m = \sum_n \mathbf{X}_{mn}\theta_n + z_m$, where $y_m$ is the price change of asset $m$, $\mathbf{X}_{mn}$ is the exposure of asset $m$ to factor $n$, $\theta_n$ are the returns of factor $n$, and $z_m$ is noise in asset $m$. We can express the linear model in matrix vector form, $\mathbf{y} = \mathbf{X}\theta + \mathbf{z}$, where $\mathbf{X}$ is an input data matrix, by assigning $y_m$ as the $m$-th entry of the vector $\mathbf{y}$, $\theta_n$ as the $n$-th entry of the vector $\theta$, and $\mathbf{X}_{mn}$ as the element of the matrix $\mathbf{X}$ in row $m$ and column $n$ The goal is to estimate $\theta$ from $\mathbf{y}$, $\mathbf{X}$, and possible statistical knowledge about $\theta$ and $\mathbf{z}$. We can see that financial prediction based on linear models relies on solving linear inverse problems. That said, some factors relate to returns in a nonlinear way [5], and financial prediction could be improved using nonlinear schemes.

Nonlinear modeling can also be used in health informatics [1], where $\mathbf{y}$ could measure patients' medical condition, $\mathbf{X}$ contains nonlinear exposure terms, and $\theta$ are explanatory variables that drive the patients' condition. Our goal is to understand the relationships between explanatory variables and patients' medical condition.

**Main idea and contributions:** In this paper, as a first step toward learning nonlinear functions, we cast them as linear inverse problems using *polynomial kernels* [6, 7] (Sec. III). Incorporating kernels into the matrix maps the nonlinear signal estimation procedure into a linear inverse problem but with an increased feature set, where the features are no longer *independent and identically distributed* (i.i.d.). Unfortunately, the kernels may create poorly conditioned matrices, and many solvers for linear inverse problems struggle with such matrices. Nonetheless, the matrices in our linear inverse problems often contain only mild linear correlations among columns, and are reasonably well conditioned.

While the polynomial kernels greatly increase the richness of the model class that captures the phenomena of interest, they also significantly increase the dimensionality of the features. For example, $N$ factors evaluated with quadratic kernels will become approximately $\frac{1}{2}N^2$ new factors. This large scale and well-conditioned linear inverse problem is well-suited to *approximate message passing* (AMP) [8, 9], an algorithmic framework for signal reconstruction that is asymptotically optimal for large scale linear inverse problems in the sense that it achieves best-possible reconstruction quality [10, 11]. Our AMP-based algorithms improve reconstruction quality of the coefficients vector $\boldsymbol{\theta}$, leading to better estimation of the nonlinear function.

Two AMP-based approaches are considered. The first follows a *Bayesian framework*, where we assume that the coefficients vector follows some known probabilistic structure. While the Bayesian framework is naive and limited in scope, our past work has shown that universal approaches that adapt to unknown statistical distributions can be integrated within solvers for linear inverse problems [12–14], thus bypassing the Bayesian limitation. The linear inverse problems resulting from our polynomial kernel expansion is solved using an AMP-based algorithm, whose Bayes optimality ensures that our function estimation procedure can succeed despite using fewer and noisier samples than other methods.

The second AMP-based approach uses *empirical Bayes* [15], where the coefficients vector is assumed to follow some parametric distribution, and in each iteration of AMP we plug maximum likelihood parameter estimates into a parametric Bayesian denoiser.

The resulting algorithms will allow data to better model dependencies between explanatory variables and phenomena of interest. These algorithms could also help reconstruct signals acquired by nonlinear analog systems, allowing hardware designers to exploit nonlinearities rather than avoid them.

**Organization:** The rest of the paper is organized as follows. Section II provides background content. Details of our approach for estimating multivariate nonlinear functions appear in Section III. Numerical results appear in Section IV, and Section V concludes.

## II. BACKGROUND

### A. Inverse problems

We present a flexible formulation for nonlinear function estimation in the form of a nonlinear *inverse problem*. We observe $M$ independent samples of the form $\{(\mathbf{x}_m, y_m)\}_{m \in \{1,...,M\}}$, where $(\mathbf{x}_m, y_m) \in \mathbb{R}^N \times \mathbb{R}$, through a nonlinear function $f(\cdot)$ and additive noise,

$$y_m = f(\mathbf{x}_m) + z_m, \tag{1}$$

for all $m \in \{1, \ldots, M\}$. In other words, the input data matrix $\mathbf{X} \in \mathbb{R}^{M \times N}$, where $\mathbf{X}$ are locations of samples, will be processed by applying a multivariate operator, $\mathbf{f}(\cdot) : \mathbb{R}^{M \times N} \to \mathbb{R}^M$, such that $\mathbf{f}$ applies $f$ on each individual row of the data matrix $\mathbf{x}$, with additive noise, $\mathbf{z} \in \mathbb{R}^M$, resulting in noisy measurements,

$$\mathbf{y} = \mathbf{f}(\mathbf{X}) + \mathbf{z} \in \mathbb{R}^M. \tag{2}$$

While the reader is likely familiar with linear inverse problems, where the operator $\mathbf{f}$ boils down to multiplication by a coefficients vector $\boldsymbol{\theta}$, i.e., $\mathbf{f}(\mathbf{X}) = \mathbf{X}\boldsymbol{\theta}$, our main interest is in nonlinear inverse problems.

We highlight that many "rules of thumb" that the sparse signal processing community has claimed, for example that sparse signals can be reconstructed from a small number of linear measurements, $M < N$, may break down when the measurement noise $\mathbf{z}$ is large or the operator $\mathbf{f}(\cdot)$ contains significant nonlinearities.

### B. Approximate message passing (AMP)

One approach for solving linear inverse problems is AMP [8, 9], which is an iterative algorithm that successively converts the matrix problem to scalar channel denoising problems with *additive white Gaussian noise* (AWGN). AMP is a fast approximation to precise message passing (cf. Baron et al. [16], Montanari [9], and references therein), and has received considerable attention because of its fast convergence and the *state evolution* (SE) formalism [8, 9, 17], which characterizes how the *mean squared error* (MSE) achieved by the next iteration of AMP can be predicted using the MSE performance of the denoiser being used. AMP solves the following linear inverse problem,

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta} + \mathbf{z}, \tag{3}$$

where the empirical *probability density function* (pdf) of $\boldsymbol{\theta}$ follows $p_{\boldsymbol{\theta}}(\boldsymbol{\theta})$, the operator $\mathbf{f}(\mathbf{X})$ multiplies $\mathbf{X}$ by the unknown coefficients vector $\boldsymbol{\theta}$, and $\mathbf{z}$ is AWGN with variance $\sigma_Z^2$. Although the AMP literature mainly considers i.i.d. Gaussian matrices, approaches such as damping [18] and Swept AMP [19] have been proposed to deal with more general matrices. After initializing $\boldsymbol{\theta}^0$ and $\mathbf{r}^0$, AMP [8, 9] proceeds iteratively according to

$$\boldsymbol{\theta}^{t+1} = \eta^t(\mathbf{X}^T \mathbf{r}^t + \boldsymbol{\theta}^t),$$
$$\mathbf{r}^t = \mathbf{y} - \mathbf{X}\boldsymbol{\theta}^t + \frac{1}{R}\mathbf{r}^{t-1}\langle\eta^{t-1'}(\mathbf{X}^T \mathbf{r}^{t-1} + \boldsymbol{\theta}^{t-1})\rangle, \tag{4}$$

where $(\cdot)^T$ denotes the transpose,

$$R = M/N$$

is the *measurement rate*, $\eta^t(\cdot)$ is a *denoising function*, and $\langle\mathbf{u}\rangle = \frac{1}{N}\sum_{i=1}^N u_i$ for some vector $\mathbf{u} \in \mathbb{R}^N$. The denoising function $\eta^t(\cdot)$ operates in a

symbol-by-symbol manner (also known as *separable*) in the original derivation of AMP [8,9]. That is, $\eta^t(\mathbf{u}) = (\eta^t(u_1), \eta^t(u_2), ..., \eta^t(u_N))$ and $\eta^{t'}(\mathbf{u}) = (\eta^{t'}(u_1), \eta^{t'}(u_2), ..., \eta^{t'}(u_N))$, where $\eta^{t'}(\cdot)$ denotes the derivative of $\eta^t(\cdot)$.

A useful property of AMP in the *large system limit* $(N, M \to \infty$ with the measurement rate $R$ constant) is that at each iteration, the vector $\mathbf{X}^T \mathbf{r}^t + \boldsymbol{\theta}^t \in \mathbb{R}^N$ in (4) is equivalent to the unknown coefficients vector $\boldsymbol{\theta}$ corrupted by AWGN. This property is based on the decoupling principle [10, 20, 21], which states that the posterior of a linear inverse problem (3) is statistically equivalent to a scalar channel. We denote the equivalent scalar channel at iteration $t$ by

$$\mathbf{q}^t = \mathbf{X}^T \mathbf{r}^t + \boldsymbol{\theta}^t = \boldsymbol{\theta} + \mathbf{v}^t, \qquad (5)$$

where $v_i^t \sim \mathcal{N}(0, \sigma_t^2)$, and $\mathcal{N}(\mu, \sigma^2)$ is a Gaussian pdf with mean $\mu$ and variance $\sigma^2$. AMP with separable denoisers, which are optimal for i.i.d. signals, has been rigorously proved to obey SE [17]. However, we will see in Section II-C that non-i.i.d. signals can be denoised better using non-separable denoisers.

Another useful property of AMP in the large system limit involves a Bayesian setting where a prior distribution for the coefficients vector $\boldsymbol{\theta}$ is available. In such Bayesian settings, AMP can use denoiser functions $\eta^t(\cdot)$ that minimize the MSE in each iteration $t$ [17]. Using such MSE-optimal denoisers, the MSE performance of AMP (4) approaches the *minimum mean squared error* (MMSE) as $t$ is increased.

### C. Non-scalar denoisers

While i.i.d. signals can be denoised in a scalar separable fashion within AMP, where each signal entry is denoised using the same scalar denoiser, real-world signals often contain dependencies between signal entries. For example, adjacent pixels in images are often similar in value, and scalar separable denoisers ignore these dependencies. Therefore, we apply non-separable denoisers to process non-i.i.d. signals within AMP. For example, if $\boldsymbol{\theta}$ is a time series containing dependencies between adjacent entries, then we can use a sliding window denoiser that processes entry $n$ of $\boldsymbol{\theta}$ using information from its neighbors [14, 22].

We will see in Section III that our signal reconstruction problem includes several types of coefficients in $\boldsymbol{\theta}$, and we expect dependencies between coefficients. Therefore, non-scalar denoisers will be used within AMP to process non-i.i.d. coefficients.

### III. Learning nonlinear functions

Having reviewed relevant background material, we now recast nonlinear inverse problems (2) as linear inverse problems (3) using polynomial kernels [6, 7],

which replace our input data matrix $\mathbf{X}$ with transformations of $\mathbf{X}$ [23].

Our nonlinear model (2) is motivated by the inadequacy of linear relationships in some applications. One example involves bioinformatics, where genetic factors involve multiplicative interactions among genes [24]. Another application involving financial prediction [2–4], where the research and development expenditures of a firm correlate with future returns in a nonlinear way [5]. Similar ideas have been widely used in the machine learning community under the context of polynomial kernel learning [6, 7], and the kernel trick has been introduced to linear inverse problems by Qi and Hughes [25]. A related model that learns interactions among variables is the multi-linear model [26], where columns that involve auto-interaction are removed from the polynomial model.

### A. Basis expansion

Recall that in our inverse problem, $\mathbf{y} = \mathbf{f}(\mathbf{X}) + \mathbf{z}$, we define measurement $m \in \{1, \ldots, M\}$ as $y_m = f(\mathbf{x}_m) + z_m$ (1). Linear inverse problems make use of models that are linear in the input factors; they are mathematically and algorithmically tractable, and can be interpreted as a first-order Taylor approximation to $f(\mathbf{x})$ [23]. However, in many applications, the true function $f(\mathbf{x})$ is far from linear in $\mathbf{x}$.

A basis function expansion replaces $\mathbf{x}$ with transformations of $\mathbf{x}$ [23]. For $\ell \in \{1, 2, \ldots, L\}$, $f(\mathbf{x})$ is expressed as in the linear basis expansion of $\mathbf{x}$:

$$f(\mathbf{x}) = \sum_{\ell=1}^{L} \theta_\ell g_\ell(\mathbf{x}).$$

This model is linear in the new variable $g_\ell(\mathbf{x})$, and $\theta_\ell$ are the coefficients. Basis expansions allow us to use a linear model to characterize and analyze nonlinear functions.

### B. Polynomial regression

We form a polynomial regression problem by applying a Taylor expansion to the multivariate nonlinear function $f(\cdot)$ [24]. In polynomial regression, we add to the original columns of the measurement matrix $\mathbf{X}_Q$, which represent individual explanatory variables, extra columns that represent interactions among variables.

Let us elaborate on the quadratic case. While we will provide details of a matrix $\mathbf{X}_Q$, that supports a quadratic Taylor expansion (6), the reader should be able to employ this concept for cubic expansions and beyond. For each

$$\mathbf{X}_Q = \begin{bmatrix} 1 & x_{11} \cdots x_{1N} & x_{11}^2 \ldots x_{1N}^2 & x_{11}x_{12} \ldots x_{1(N-1)}x_{1N} \\ 1 & x_{21} \ldots x_{2N} & x_{21}^2 \ldots x_{2N}^2 & x_{21}x_{22} \ldots x_{2(N-1)}x_{2N} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_{M1} \ldots x_{MN} & x_{M1}^2 \ldots x_{MN}^2 & x_{M1}x_{M2} \ldots x_{M(N-1)}x_{MN} \end{bmatrix}. \tag{6}$$

measurement, we use a Taylor expansion of the $N$ factor variables:

$$
\begin{aligned}
y_m = {} & \theta_1 + \sum_{n=1}^{N}[\boldsymbol{\theta_2}]_n x_{mn} + \sum_{n=1}^{N}[\boldsymbol{\theta_3}]_n x_{mn}^2 \\
& + \sum_{n_1=1}^{N}\sum_{n_2=n_1+1}^{N}[\boldsymbol{\theta_4}]_n x_{mn_1} x_{mn_2},
\end{aligned}
\tag{7}
$$

where $\theta_1$ is a constant, $\boldsymbol{\theta_2}, \boldsymbol{\theta_3} \in \mathbb{R}^N$ are coefficient vectors for linear and quadratic terms, respectively, $\boldsymbol{\theta_4} \in \mathbb{R}^{\frac{N(N-1)}{2}}$ is a coefficient vector for cross terms, and the subscript $n$ in $[\boldsymbol{\theta_4}]_n$ depends on $n_1$ and $n_2$.

Our quadratic Taylor approximation is a basis expansion, where we have chosen $g(\mathbf{x})$ as follows: (*i*) $g(\mathbf{x}) = 1$ corresponds to a DC constant (*ii*) $N$ linear terms corresponding to the original data, $g(\mathbf{x}) = x_n$, $n \in \{1, \ldots, N\}$; (*iii*) $N$ quadratic terms corresponding to squares of individual linear terms, $g(\mathbf{x}) = (x_n)^2$; and (*iii*) $\frac{N(N-1)}{2}$ cross terms corresponding to products of pairs of linear terms, $g(\mathbf{x}) = x_{n_1} x_{n_2}$, where $n_2 > n_1$, $n_1, n_2 \in \{1, \ldots, N\}$. We assume that the features matrix, $\mathbf{X}$, is i.i.d. zero mean Gaussian for ease of analysis; different types of $\mathbf{X}$ are left for future work.

The polynomial regression model is formulated as a linear inverse problem (3) in matrix vector form,

$$\mathbf{y} = \mathbf{X}_Q\boldsymbol{\theta} + \mathbf{z} = \mathbf{X}_Q \begin{bmatrix} \theta_1 \\ \boldsymbol{\theta_2} \\ \boldsymbol{\theta_3} \\ \boldsymbol{\theta_4} \end{bmatrix} + \mathbf{z},$$

where $\boldsymbol{\theta} \in \mathbb{R}^L$ is the coefficient vector, and $L$ is evaluated below (8). In our matrix $\mathbf{X}_Q$ (6), each row is an instance or sample, and each column is an attribute or feature.

Our goal is to estimate the regression coefficients in the vector $\boldsymbol{\theta}$ from $\mathbf{X}_Q$ and $\mathbf{y}$. The measurement matrix $\mathbf{X}_Q \in \mathbb{R}^{M \times L}$ will include one DC column, $N$ linear term columns, $N$ quadratics (squared column), and $\frac{N(N-1)}{2}$ cross terms. This matrix has the form (6), and it can be seen that

$$L = 1 + 2N + \frac{N(N-1)}{2}. \tag{8}$$

To solve this linear inverse problem using an AMP-based approach, we normalize each column of $\mathbf{X}_Q$,

$[\mathbf{X}_Q]_\ell$ to have unit norm, where $\ell \in \{1, 2, \ldots, L\}$, and denote this normalized matrix by $\mathbf{X}'_Q$,

$$\mathbf{y} = \mathbf{X}_Q\boldsymbol{\theta} + \mathbf{z} = \mathbf{X}'_Q\boldsymbol{\theta}' + \mathbf{z}, \tag{9}$$

where each entry of $[\mathbf{X}_Q]_\ell$ obeys

$$[\mathbf{X}'_Q]_{\ell m} = \frac{[\mathbf{X}_Q]_{\ell m}}{||[\mathbf{X}_Q]_\ell||_2},$$

and the regression coefficients satisfy

$$\theta'_\ell = \theta_\ell ||[\mathbf{X}_Q]_\ell||_2. \tag{10}$$

### C. SVD of normalized quadratic matrix $\mathbf{X}'_Q$

While the normalized matrix $\mathbf{X}'_Q$ converts our quadratic nonlinear inverse problem into a linear one, it contains dependencies between linear and quadratic columns as well as between the linear and cross terms. Unfortunately, it is well known that many solvers for linear inverse problems struggle with such matrices.

Surprisingly, our matrix (6) works well within some AMP-based approaches, as will be demonstrated by numerical results in Section IV. Why does our matrix perform well within AMP? *Despite containing dependencies between columns, these dependencies are nonlinear in nature, and linear correlations between columns turn out to be mild*. In fact, a *singular value decomposition* (SVD) of $\mathbf{X}'_Q$ reveals that it is reasonably well-conditioned. In particular, we have seen numerically that most of the *singular values* (SVs) seem to follow the semicircle law. That said, the first (largest) SV is larger than suggested by the semicircle law.

To see why the first SV, $\sigma_1$, is larger, recall that $\mathbf{X}'_Q$ is comprised of one DC column, $N$ linear term columns, $N$ quadratic ones, and $\frac{N(N-1)}{2}$ cross term columns. Because $\mathbf{X}'_Q$ has unit norm columns, entries of the DC column are $1/\sqrt{M}$, and so the sum of elements of the first column is $\sqrt{M}$. The $N$ quadratic columns are non-negative, and because they too have unit norm, the average squared value is $1/M$, suggesting that the average is $\Theta(1/\sqrt{M})$. The sums of elements of all $N$ linear and $\frac{N(N-1)}{2}$ cross term columns are near zero, because these are zero mean Gaussian *random variables* (RVs), and products of zero mean Gaussian RVs, respectively. We see that the first SV, $\sigma_1$, corresponds to an all constant (or roughly all constant) column multiplied by a row that contains significant non-zero entries corresponding to the DC column and $N$ quadratic columns, while row

entries corresponding to linear and cross term columns are close to zero.

Under some assumptions, we can estimate the amount of energy represented by the first SV, $\sigma_1^2$. Suppose that the original linear columns are Gaussian, $X \sim \mathcal{N}(0,1)$. Under this assumption, the quadratic element $\chi = X^2$ has a chi-squared distribution, where $E[\chi] = E[X^2] = 1$ and $\text{var}[\chi] = 2$. Therefore, $E[\chi^2] = E[\chi]^2 + \text{var}(\chi) = 3$. As we will need to normalize individual entries of quadratic terms by roughly $\sqrt{3M}$, the average energy of the DC component of these columns is $1/3$. Similarly, it can be shown that linear and cross term columns have average energy $1/M$ aligned with the first singular column vector. In summary, the energy in $\sigma_1^2$ is comprised of (*i*) unit energy for the DC column; (*ii*) $N/M$ for the $N$ linear columns; (*iii*) $N/3$ for the $N$ quadratic ones; and (*iv*) $\frac{N(N-1)}{2M}$ for cross term columns. Therefore, we predict the total energy in $\sigma_1$ to obey

$$\sigma_{1,pred}^2 = 1 + N/3 + \frac{N(N+1)}{2M}. \qquad (11)$$

Our analysis of the first singular value is inaccurate, because the first singular vector column is only roughly constant, and while computing the SVD this column is modified in order to maximize the energy of the first rank-one component. Therefore, $\sigma_{1,pred}^2$ can be interpreted as a *lower bound* for $\sigma_1^2$. That said, numerical experiments presented in Table I show that our prediction (11) provides a reasonable approximation. In the table, results for several $(M, N)$ pairs are provided. For each pair, we average empirical values for $\sigma_1^2$, the energy in the first SV, over 20 matrices; these empirical averages are compared to the prediction (11). It can be seen that $\sigma_1^2$ is typically larger by 0.6–0.75; seeing that unit norm columns in the normalized matrix $\mathbf{X}'_Q$ imply that the average SV has unit energy, this extra energy seems plausible.

Finally, although we have focused on the normalized quadratic matrix, $\mathbf{X}'_Q$, in further numerical work (not reported here) we evaluated a cubic matrix with normalized columns. It too has an SVD where $\sigma_1$ is larger while other SVs seem to follow the semicircle law.

### D. AMP-based algorithm

We solve our linear inverse problem (9) using AMP, where two points should be highlighted. First, our denoiser can incorporate the Bayesian prior information. Specifically, we use conditional expectation denoisers that minimize the MSE [17]. Second, owing to the structure of our matrix (Section III-C), various AMP variants that promote convergence can be used [18, 19, 27]. That said, these variants all have their shortcomings, and possible divergence of AMP should be tracked carefully.

TABLE I
EMPIRICAL VALUE OF $\sigma_1^2$ COMPARED TO OUR PREDICTION (11).

| $M$ | $N$ | $L$ (8) | $\sigma_1^2$ (empirical) | $\sigma_{1,pred}^2$ (11) |
|---|---|---|---|---|
| 1000 | 10 | 66 | 4.99 | 4.39 |
| 1500 | 15 | 136 | 6.72 | 6.08 |
| 2000 | 20 | 231 | 8.41 | 7.77 |
| 3000 | 20 | 231 | 8.35 | 7.74 |
| 3000 | 30 | 496 | 11.84 | 11.16 |
| 4000 | 40 | 861 | 15.23 | 14.54 |
| 4500 | 50 | 1326 | 18.63 | 17.95 |
| 5000 | 60 | 1891 | 22.09 | 21.37 |
| 5500 | 70 | 2556 | 25.51 | 24.79 |
| 5000 | 80 | 3321 | 29.06 | 28.31 |
| 6000 | 80 | 3321 | 28.94 | 28.21 |
| 8000 | 80 | 3321 | 28.78 | 28.07 |
| 8000 | 90 | 4486 | 32.26 | 31.51 |
| 6000 | 100 | 5151 | 35.93 | 35.18 |
| 8000 | 100 | 5151 | 35.66 | 34.96 |

## IV. NUMERICAL RESULTS

Our construction of the quadratic polynomial regression model in Section III results in a linear inverse problem (3) whose solution forms an estimate of a multivariate nonlinear function (2) that relates the explanatory variables to the phenomena of interest. This resulting linear inverse problem will now be solved by two AMP-based approaches, Bayesian AMP and empirical Bayes.

### A. Bayesian AMP

**Non-i.i.d. model for $\theta$:** Our Bayesian approach considers four groups of coefficients (7), where $\theta_1 \in \mathbb{R}$, $\theta_2, \theta_3 \in \mathbb{R}^N$, and $\theta_4 \in \mathbb{R}^{\frac{N(N-1)}{2}}$ are the DC, linear, quadratic, and cross term coefficients, respectively. We modeled each individual entry among these $L$ coefficients as *Bernoulli Gaussian* (BG), where the Bernoulli part is a probability $p$ that the entry is nonzero, in which case its distribution is zero mean Gaussian with some variance. To be specific, (*i*) our DC coefficent obeys $\theta_1 \sim \mathcal{N}(0, 10)$, meaning that it is zero mean Gaussian with variance 10; (*ii*) each entry among the $N$ linear term coefficients satisfies $[\theta_2]_n \sim 0.2\mathcal{N}(0,1) + 0.8\delta_0$, i.e., zero mean unit norm Gaussian with probability 0.2, else zero; (*iii*) the $N$ quadratic term coefficients obey $[\theta_3]_n \sim 0.2\mathcal{N}(0, 0.5) + 0.8\delta_0$; and (*iv*) for the $\frac{1}{2}N(N-1)$ cross term coefficients, $[\theta_4]_n \sim 0.03\mathcal{N}(0, 0.1) + 0.97\delta_0$. Although the four groups of coefficients have different distributions, all $L$ entries that follow this model are statistically independent.

**Baseline LASSO algorithm:** The baseline algorithm used to solve (9) is the *least absolute shrinkage and selection operator* (LASSO) [28], which minimizes the sum of squared errors subject to a constraint on the $\ell_1$ norm of the coefficients [28]. In our polynomial model,

the LASSO estimator $\widehat{\boldsymbol{\theta}}$ is calculated in Lagrangian form:

$$\widehat{\boldsymbol{\theta}} = \frac{1}{2}\operatorname*{argmin}_{\boldsymbol{\theta}} ||\mathbf{y} - \mathbf{X}_Q\boldsymbol{\theta}||_2^2 + \sum_{j=1}^{4}\lambda_j\|\boldsymbol{\theta}_j\|_1, \quad (12)$$

where $\lambda_1, \ldots, \lambda_4$ are tuning parameters. In principle, we could perform grid search over all four parameters, $\lambda_1, \ldots, \lambda_4$, but it is computationally intractable. Therefore, we report the performance obtained by setting all parameters to be equal, which reduces the search space.

**AMP-based approach:** As a proof of concept, we have designed a denoiser specifically for our non-i.i.d. model. Because all $L$ entries that follow this model are statistically independent, we used $L$ scalar denoisers. However, because individual entries among our four groups of coefficients, $\theta_1 \in \mathbb{R}$, $\boldsymbol{\theta_2}, \boldsymbol{\theta_3} \in \mathbb{R}^N$, and $\boldsymbol{\theta_4} \in \mathbb{R}^{\frac{N(N-1)}{2}}$ follow different distributions, four different scalar denoisers were used. Details of Bayesian denoisers for BG signals appear in [29].

**Signal generation:** We evaluate the performance of AMP in the Bayesian setting, which is a planted inference problem. The experiment allows us to validate the suitability of AMP for the quadratic basis, e.g. (6).

We generated the feature matrix, $\mathbf{X}$, as i.i.d. Gaussian with dimension $N = 100$. These linear terms were then transformed into a quadratic form $\mathbf{X}_Q'$ with normalized columns (9). The number of columns in the normalized matrix was $L = 5151$ (8), and the number of rows $M = 5400$, Next, we created quadratic multivariate functions by generating $\boldsymbol{\theta}$ vectors following our non-i.i.d. model. The expected energy of each group of coefficients satisfies $E_{DC} = 10$, $E_{linear} = 0.2 \times N = 20$, $E_{quadratic} = 0.2 \times N \times 0.5 = 10$, and $E_{cross} = 0.03 \times \frac{N^2 - N}{2} \times 0.1 = 14.85$. Finally, the measurement noise $\mathbf{z}$ was AWGN with variance $\sigma_Z^2 = 0.004$.

**MSE performance:** Fig. 1 shows the MSE performance for estimated coefficients, $\boldsymbol{\theta}$. We estimated the coefficients using LASSO, *swept AMP* (SwAMP) [19] and *vector AMP* (VAMP) [27]. The left panel of the figure shows the MSE obtained when estimating the original coefficients $\boldsymbol{\theta}$, where the estimator $\widehat{\theta}$ can be calculated using (13),

$$\widehat{\theta}_\ell = \frac{\widehat{\theta'}_\ell}{||[\mathbf{X}_Q]_\ell||_2}, \quad (13)$$

$l \in \{1, \ldots, L\}$, and $\widehat{\boldsymbol{\theta}'}$ are estimated coefficients of $\boldsymbol{\theta}'$. SwAMP and VAMP both converge well for normalized quadratic matrices. However, it can be seen in Fig. 1 that VAMP requires less than one hundred iterations to converge; SwAMP requires a few hundred, and its individual iterations require more computation than those of VAMP; our specific implementation of LASSO requires thousands of iterations. Because our AMP based approaches are expected to be Bayes optimal

while LASSO does not share these optimality properties, there is no surprise that AMP-based approaches obtain lower MSE.

To make sure that our function reflects the nonlinear function well, the right panel of Fig. 1 shows the MSE obtained when applying our estimated polynomial function to predict test data,

$$\frac{||\mathbf{y}_{test} - \mathbf{X}_{test}\widehat{\boldsymbol{\theta}}||_2^2}{K} = \frac{||\mathbf{X}_{test}(\boldsymbol{\theta} - \widehat{\boldsymbol{\theta}})||_2^2}{K},$$

where we held back $K = 600$ test measurements (recall that $M = 5400$), $\mathbf{X}_{test} \in \mathbb{R}^{K \times L}$ has the same format as $\mathbf{X}_Q$, and $\mathbf{y}_{test} \in \mathbb{R}^K$. Note that the MSE for coefficients, $\boldsymbol{\theta}$, is inapplicable to real-world problems, because the true coefficients do not exist, and we are merely modeling some nonlinear dependence as a low-order Taylor series. In our synthetic experiment, we are using the MSE over the test data as a metric of interest.

### B. Empirical Bayes

**Nonlinear function:** Nonlinear function learning is now performed using empirical Bayes within AMP [15]. We employ the quadratic formulation (9) and learn the coefficients vector $\boldsymbol{\theta}$ to approximate a family of (mildly) nonlinear functions,

$$\boldsymbol{y} = \sum_{i=1}^{3} w_i \sin\left(\boldsymbol{X}\boldsymbol{\rho}_i + \boldsymbol{\phi}_i\right) + \mathbf{z}, \quad (14)$$

where $w_1 = 0.1$, $w_2 = 0.3$, and $w_3 = 0.6$ are weights of the sinusoids, $\boldsymbol{\rho}_i \in \mathbb{R}^N$ is a BG vector, $\boldsymbol{X}\boldsymbol{\rho}_i \in \mathbb{R}^M$, $\boldsymbol{\phi} \in \mathbb{R}^M$ are phase shifts uniformly distributed between 0 and $2\pi$, the sine is applied element-by-element, and the noise $\boldsymbol{z} \in \mathbb{R}^M$ is AWGN with variance $10^{-4}$. Note that the vectors $\boldsymbol{\rho}_i$ are chosen to be sparse BG, in order for the coefficients vector $\boldsymbol{\theta}$ fit by AMP to the quadratic expansion to also be sparse.

**AMP-based empirical Bayes:** In contrast to the Bayesian case, we assume that $\boldsymbol{\theta}_2$, $\boldsymbol{\theta}_3$, and $\boldsymbol{\theta}_4$ are BG, and their parameters are estimated using *maximum likelihood* (ML) in each AMP iteration. The DC coefficient $\theta_1$ is assumed to be Gaussian. The ML parameters are plugged into Bayesian denoisers for the 4 components.

**MSE performance:** We generated nonlinear functions and ran our empirical Bayes algorithm, LASSO, and a pseudoinverse approach (least squares). Each run of LASSO requires many iterations, and we use cross validation to regularize the parameter selection procedure. AMP with damping requires fewer iterations than LASSO. Empirical results for different measurement rates, $R = M/L$, appear in Table. II. AMP obtains lower MSE than LASSO, which in turn obtains lower MSE than pseudoinverse.
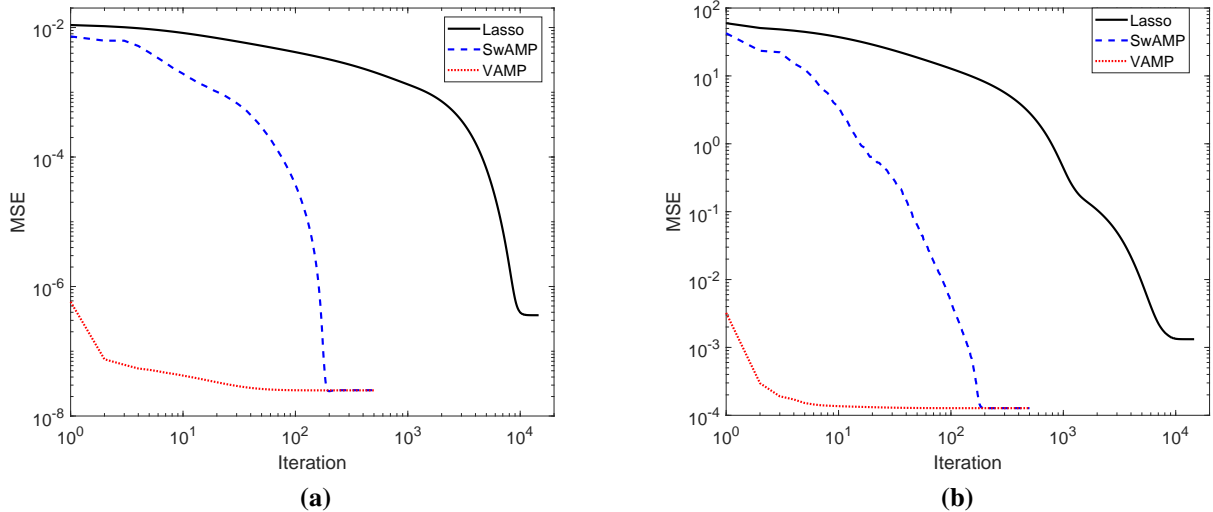
Fig. 1. Performance of LASSO, SwAMP, and VAMP. The MSE is shown in the vertical axis, while the horizontal axis reflects the iteration number, $t$. Left panel **(a)**: MSE performance in recovering the unknown coefficients, $\boldsymbol{\theta}$. Right panel **(b)**: MSE performance in predicting the test data.

TABLE II
EMPIRICAL MSE ON TEST DATA FOR NONLINEAR FUNCTION
ESTIMATION.

| Measurement Rate $R = \frac{M}{L}$ | Median MSE over 20 Realizations | | |
|---|---|---|---|
| | **LASSO** | **AMP** | **Pseudoinverse** |
| 0.14 | 0.0382 | 0.0293 | 0.041 |
| 0.28 | 0.0298 | 0.0228 | 0.033 |
| 0.56 | 0.0063 | 0.0036 | 0.01 |

## V. DISCUSSION

In this paper, we studied nonlinear function estimation, where a nonlinear function of interest is regressed on a set of features. We linearized the problem by considering low-order polynomial kernel expansion, and solved the resulting linear inverse problem using *approximate message passing* (AMP). Numerical results confirm that our AMP-based approaches learn the function better than the widely used *least absolute shrinkage and selection operator* (LASSO) [28], offering markedly lower error in predicting test data for both Bayesian and non-Bayesian settings.

While we have presented a first step toward estimating nonlinear functions by appling AMP to polynomial regression, many open problems remain.

**Dependencies between coefficients:** In past work, we used non-scalar sliding window denoisers to process coefficient vectors $\boldsymbol{\theta}$ that contained dependencies between entries [14, 22]. It is not clear whether similar dependencies will appear in our $\boldsymbol{\theta}$. While it seems plausible that exposure weights corresponding to the $N$ original

columns, the $N$ quadratic terms, and $N(N-1)/2$ cross terms will have different distributions, it is not clear whether each group is i.i.d. or contains intra-group dependencies. In ongoing work, we are processing all terms corresponding to the same original column (the original column, its quadratic, and $N-1$ associated product columns) together, which could be processed with block denoising. This form of joint processing will support possible dependencies between lower order Taylor coefficients and higher order ones; such dependencies have been noted between parent and children wavelet coefficients [30].

**Other kernels:** In this paper, we considered a second-order polynomial kernel. Future work will naturally extend to selecting the degree of the polynomial kernel as well. Further, we will consider other widely used kernels.

**Results on real datasets:** While we reported promising results for nonlinear function estimation with AMP in Bayesian and empirical Bayes settings, the performance of our algorithms must be tested on real datasets. In these datasets, various problems may appear, for example the prior is unavailable; the measurement matrix may be poorly conditioned; the function of interest may not belong to the hypothesis class; and the noise may be heavy tailed [2], resulting in a mismatched estimation problem. We will explore the application of more advanced adaptive variants of AMP in the absence of a known prior [12–14]. When the true function does not belong to the hypothesis class, which are polynomials of degree two or three in this paper, the best one can hope for is to recover the function of interest up to a

projection error onto the hypothesis class. We will also explore the usual bias/variance trade-offs that arise in such settings.

**Nonlinear acquisition and reconstruction:** Since the work of Gauss and his contemporaries [31], hardware designers have been keenly aware that the mathematics involved in processing linearly obtained measurements is more mature than that for nonlinear measurements. However, algorithms that estimate multivariate nonlinear functions can also be used to reconstruct signals measured nonlinearly. The same polynomial kernels [6, 7] used above to expand the matrix can also be used to approximate a nonlinear function with a linear one. Such advances will allow designers to stop worrying about the nonlinearities inherent in many hardware systems.

## REFERENCES

[1] L. A. Dalton and E. R. Dougherty, "Optimal classifier with minimum expected error within a Bayesian framework - Part 1: Discrete and Gaussian model," *Pattern Recognition*, vol. 46, pp. 1301–1314, Nov. 2012.

[2] R. C. Grinold and R. N. Kahn, *Active portfolio management: a quantitative approach for providing superior returns and controlling risk*. McGraw-Hill Companies, 2000.

[3] E. Fama and K. French, "Common risk factors in the returns on stocks and bonds," *J. Finan. Econ.*, vol. 33, no. 1, pp. 3–56, 1993.

[4] N. Jegadeesh and S. Titman, "Returns to buying winners and selling losers: Implications for stock market efficiency," *J. Finance*, vol. 48, no. 1, pp. 65–91, 1993.

[5] L. K. Chan, J. Lakonishok, and T. Sougiannis, "The stock market valuation of research and development expenditures," National Bureau of Economic Research, Tech. Rep., 1999.

[6] J. Fan, N. E. Heckman, and M. P. Wand, "Local polynomial kernel regression for generalized linear models and quasi-likelihood functions," *J. Amer. Stat. Assoc.*, vol. 90, no. 429, pp. 141–150, 1995.

[7] K. I. Kim, K. Jung, and H. J. Kim, "Face recognition using kernel principal component analysis," *IEEE Signal Process. Lett.*, vol. 9, no. 2, pp. 40–42, 2002.

[8] D. L. Donoho, A. Maleki, and A. Montanari, "Message passing algorithms for compressed sensing," *Proc. Nat. Academy Sci.*, vol. 106, no. 45, pp. 18 914–18 919, Nov. 2009.

[9] A. Montanari, "Graphical models concepts in compressed sensing," *Compressed Sensing: Theory and Applications*, pp. 394–438, 2012.

[10] D. Guo, D. Baron, and S. Shamai, "A single-letter characterization of optimal noisy compressed sensing," in *Proc. Allerton Conf. Commun., Control, and Comput.*, Sept. 2009, pp. 52–59.

[11] S. Rangan, A. K. Fletcher, and V. K. Goyal, "Asymptotic analysis of MAP estimation via the replica method and applications to compressed sensing," *CoRR*, vol. abs/0906.3234, June 2009.

[12] D. Baron and M. F. Duarte, "Universal MAP estimation in compressed sensing," in *Proc. Allerton Conf. Commun., Control, and Comput.*, Sept. 2011, pp. 768–775.

[13] J. Zhu, D. Baron, and M. F. Duarte, "Recovery from linear measurements with complexity-matching universal signal estimation," *IEEE Trans. Signal Process.*, vol. 63, no. 6, pp. 1512–1527, Mar. 2015.

[14] Y. Ma, J. Zhu, and D. Baron, "Approximate message passing algorithm with universal denoising and Gaussian mixture learning," *IEEE Trans. Signal Process.*, vol. 65, no. 21, pp. 5611–5622, Nov. 2016.

[15] Y. Ma, J. Tan, N. Krishnan, and D. Baron, "Empirical Bayes and full Bayes for signal estimation," *Arxiv preprint arxiv:1405.2113v1*, May 2014.

[16] D. Baron, S. Sarvotham, and R. G. Baraniuk, "Bayesian compressive sensing via belief propagation," *IEEE Trans. Signal Process.*, vol. 58, no. 1, pp. 269–280, Jan. 2010.

[17] M. Bayati and A. Montanari, "The dynamics of message passing on dense graphs, with applications to compressed sensing," *IEEE Trans. Inf. Theory*, vol. 57, no. 2, pp. 764–785, Feb. 2011.

[18] S. Rangan, P. Schniter, and A. Fletcher, "On the convergence of approximate message passing with arbitrary matrices," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, July 2014, pp. 236–240.

[19] A. Manoel, F. Krzakala, E. W. Tramel, and L. Zdeborová, "Swept approximate message passing for sparse estimation," in *Proc. Int. Conf. Machine Learning*, vol. 37, July 2015, pp. 1123–1132.

[20] D. Guo and S. Verdú, "Randomly spread CDMA: Asymptotics via statistical physics," *IEEE Trans. Inf. Theory*, vol. 51, no. 6, pp. 1983–2010, June 2005.

[21] D. Guo and C. C. Wang, "Multiuser detection of sparsely spread CDMA," *IEEE J. Sel. Areas Commun.*, vol. 26, no. 3, pp. 421–431, Apr. 2008.

[22] K. Sivaramakrishnan and T. Weissman, "A context quantization approach to universal denoising," *IEEE Trans. Signal Process.*, vol. 57, no. 6, pp. 2110–2129, June 2009.

[23] T. Hastie, R. Tibshirani, and J. H. Friedman, *The Elements of Statistical Learning*. Springer, Aug. 2001.

[24] V. Kekatos and G. Giannakis, "Sparse Volterra and polynomial regression models: Recoverability and estimation," *IEEE Trans. Signal Process.*, vol. 59, no. 12, pp. 5907–5920, Dec. 2011.

[25] H. Qi and S. Hughes, "Using the kernel trick in compressive sensing: Accurate signal recovery from fewer measurements," in *IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*. IEEE, 2011, pp. 3940–3943.

[26] B. Nazer and R. Nowak, "Sparse interactions: Identifying high-dimensional multilinear systems via compressed sensing," in *Proc. Allerton Conference Commun., Control, and Comput.*, Sept. 2010, pp. 1589–1596.

[27] S. Rangan, P. Schniter, and A. K. Fletcher, "Vector approximate message passing," in *Proc. Int. Symp. Inf. Theory (ISIT)*, July 2017, pp. 1588–1592.

[28] R. Tibshirani, "Regression shrinkage and selection via the LASSO," *J. Royal Stat. Soc. Series B (Methodological)*, vol. 58, no. 1, pp. 267–288, 1996.

[29] J. Vila and P. Schniter, "Expectation-maximization Bernoulli-Gaussian approximate message passing," in *Proc. IEEE 45th Asilomar Conf. Signals, Syst., and Comput.*, Nov. 2011, pp. 799–803.

[30] E. P. Simoncelli and E. H. Adelson, "Noise removal via Bayesian wavelet coring," in *Proc. Int. Conf. Image Process.*, vol. 1. IEEE, Sept. 1996, pp. 379–382.

[31] C. F. Gauss, *Theoria motus corporum coelestium in sectionibus conicis solem ambientium*, 1809.