

The problem: ML models hate extremes! (They're biased!)

Why Machine Learning Models Systematically Underestimate Extreme Values II: How to Fix It with LatentNN

YUAN-SEN TING^{1,2,3}

The solution: LatentNN to the rescue!

¹Department of Astronomy, The Ohio State University, Columbus, OH 43210, USA

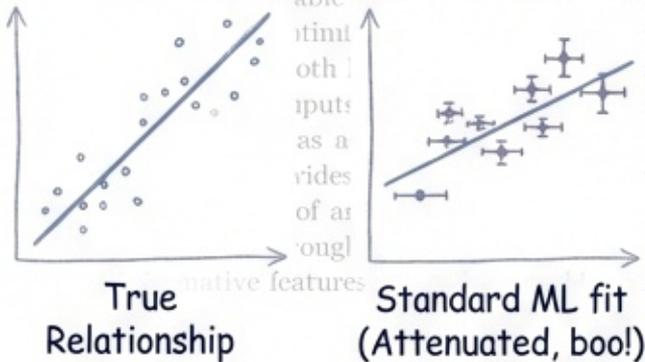
²Center for Cosmology and AstroParticle Physics (CCAPP), The Ohio State University, Columbus, OH 43210, USA

³Maz-Planck-Institut für Astronomie, Königstuhl 17, D-69117 Heidelberg, Germany

ABSTRACT

Attenuation bias—the systematic underestimation of regression errors in input variables—affects astronomical data-driven models. For linear regression, this problem was solved by treating the true input values as latent variables to be estimated alongside model parameters. In this paper, we treat the true input values as latent variables to be estimated alongside model parameters. We introduce LatentNN, a method that treats the true input values as latent variables to be estimated alongside model parameters. We demonstrate that LatentNN says: "Hold on, those noisy inputs aren't the real story. Let me figure out what the true values probably were while I'm learning!"

Basically, noisy inputs make models lazy and squish their predictions towards the average. This is ATTENUATION BIAS.



LatentNN says: "Hold on, those noisy inputs aren't the real story. Let me figure out what the true values probably were while I'm learning!"

Perfect for noisy data, like a lot of astronomy!

Keywords: methods: statistical — methods: data analysis — techniques: spectroscopy — stars: fundamental parameters

1. INTRODUCTION

Modern astronomy increasingly relies on neural networks for inference (Cranmer et al. 2020; Huertas-Company & Lanusse 2023; Ting 2025a). Data-driven approaches now routinely map observed quantities to physical parameters across a wide range of applications, from stellar parameters derived from spectra (e.g., Ness et al. 2015; Fabbro et al. 2015; Lenng & Bovy 2019; Ting et al. 2019) to photometric redshifts from imaging (Hoyle 2016; Li et al. 2022a; Lin et al. 2022; Zhou et al. 2022) to stellar properties from time series light curves (Pasquet et al. 2019; Hon et al. 2021).

With the rise of neural network approaches, however, input uncertainties are often ignored or left unmodeled. This is a design choice rather than an oversight. Classification and regression with neural networks, by contrast, arose in domains like computer vision (Simonyan & Zisserman 2014; Szegedy et al. 2014; He et al. 2016) and natural language processing (Devlin et al. 2018; Brown et al. 2020; OpenAI et al.

2023; Touvron et al. 2023) where input noise is negligible, and its standard architectures have little mechanism for direct uncertainty quantification. When applied to astronomical data, this can lead to misdiagnosis (Kelly et al. 2022; Ting 2025b). Astronomy operates in a regime that differs from most machine learning applications. Astronomical observations often have signal-to-noise ratios (SNR) of order 10, yet we seek inference precision at the percent level or better. This tension between noisy inputs and stringent precision requirements creates a unique challenge for neural network inference.

This phenomenon is known as attenuation bias (Spearman 1904; Frost & Thompson 2002; Kelly 2007). In Paper I of this series (Ting 2025b), we demonstrated this effect in detail for linear regression and explored its implications for astronomical data analysis. Attenuation bias manifests as a compression of the predicted dynamic range: high values are consistently predicted too low, and low values are consistently predicted too high. The effect depends on the ratio of input measurement error to the true value.

Astronomy has BIG data now, so we need ML to make sense of it all.



black box?

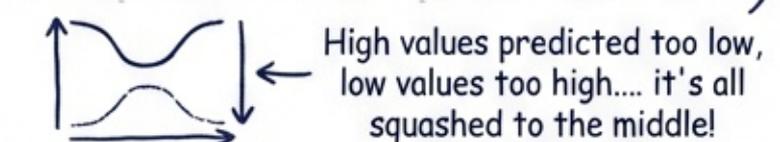
Most ML models are like:

La la la, I can't see your error bars!

?

?

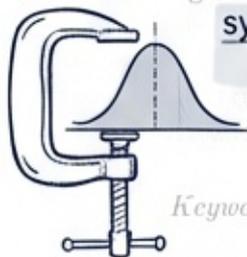
?



Especially bad when data is noisy (low SNR)!

to intrinsic signal range, and it is independent of sample size. As a rule of thumb, attenuation bias becomes important at input SNR of order 10, especially at low-dimensional problems. The exact magnitude depends on the dimensionality and correlation structure of the inputs, but the bias is never zero when input uncertainties are present. We refer interested readers to Paper I for detailed derivations and further discussion of attenuation bias in the context of linear regression.

The consequences of attenuation bias are not merely increased scatter in predictions. When we train a regression model to predict some quantity from noisy inputs, treating the observed inputs as exact is a reasonable approximation only when measurement errors are small compared to the intrinsic variation in the data. When measurement errors constitute even a modest fraction of the signal, a systematic compression of the predicted dynamic range emerges—high values are underestimated and low values are overestimated—a bias that cannot be removed by collecting more data or by using more sophisticated models. It is a fundamental consequence of treating noisy measurements as if they were exact.



systematic compression of the predicted dynamic range.

This is the 'squashing' effect we saw on the last page!

Keywords: methods: statistical fundamental parameters

1. INTRODUCTION

Modern astronomy increasingly relies on neural networks for inference (Cranmer et al. 2020; Huertas-Company & Lamuse 2023; Ting 2023a). Data-driven approaches now routinely map observed quantities to physical parameters across a wide range of applications, treat the true input values as latent variables. (g., Ness et al. 2015; Fabbro et al. 2019)

The clever trick from old-school stats: pretend the true values are hidden variables we can solve for. (Hoyle 2016; Li et al. 2022b; Pasquet et al. 2019; Hon et al. 2022)

With the rise of neural network approaches, however, extends this idea to neural networks.

And now we're bringing this powerful idea to modern deep learning!

2. ATTENUATION BIAS

To understand attenuation bias, consider a simple regression problem where we want to learn the relation-

The **TRUE** relationship we want to find.

The **BIASED** fit from standard ML.

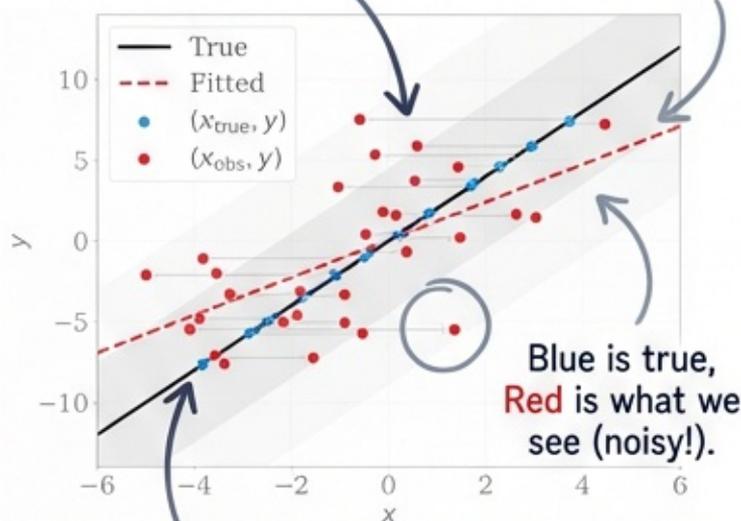


Figure 1. Schematic showing true data points (blue circles show true points, black solid line); red squares show observed data points, adding measurement error. Horizontal dashed lines indicate horizontal displacements between true and observed positions. The tilted grey shaded band represents the true relationship and $\pm 2\sigma_x$ (lighter) uncertainty regions around the true relationship. A linear fit to the observed data (red dashed line) yields a shallower slope, demonstrating how input measurement errors systematically bias regression results.

See? The fitted line is too shallow! It's been **SQUASHED** towards the middle.

This is **ATTENUATION BIAS** in action. The noise makes the model underestimate the true slope.

Basically, noise makes the data spread out more horizontally.

$$\text{Var}(x_{\text{obs}}) = \sigma_{\text{range}}^2 + \sigma_{\epsilon}^2. \quad (1)$$

This stretches the data horizontally without a corresponding vertical stretch, because the measurement errors in x do not affect y . The observed distribution appears wider in the x -direction than the true distribution, while the y -extent remains unchanged. Any regression method treating x_{obs} as exact will therefore fit a slope that is systematically shallower than the true relationship. This weakening of the apparent relationship is why attenuation bias is also known as regression dilution.

regression dilution

A.k.a. "watering down" the relationship.

This effect can be visualized by plotting data points and their regression line. Figure 1 illustrates this with a simple example: true data points (blue) following a linear relationship $y = 2x$ are scattered horizontally by measurement error to their observed positions (red), with arrows indicating the displacement. Grey tilted

THE "SQUASHING" FACTOR (MATH EDITION)

This number tells us how much the line gets flattened (always < 1)

$$\mathbb{E}[\hat{\beta}] = \beta \cdot \lambda_{\beta}, \quad \text{where} \quad \lambda_{\beta} = \frac{1}{1 + (\sigma_x/\sigma_{\text{range}})^2} \quad (2)$$

The quantity λ_{β} is called the attenuation factor. It depends only on the ratio of measurement error to the intrinsic signal range, $\sigma_x/\sigma_{\text{range}}$. When measurement errors are small compared to the signal range, λ_{β} is close to unity and the bias is negligible. When measurement errors are comparable to the signal range, λ_{β} can be much less than unity, indicating large bias. $\lambda_{\beta} < 1$ are predicted too high. In the extreme case where measurement errors dominate ($\sigma_x \gg \sigma_{\text{range}}$), λ_{β} approaches zero and the model predicts the mean for all inputs, obscuring the true relationship.

Consider some representative cases in the one-dimensional setting:

- At $\sigma_x/\sigma_{\text{range}} = 0.1$ (SNR = 10), the attenuation factor is $\lambda_{\beta} \approx 0.99$, corresponding to 1% bias (Not bad)
- At $\sigma_x/\sigma_{\text{range}} = 0.33$ (SNR ≈ 3), the factor is $\lambda_{\beta} \approx 0.90$, corresponding to 10% bias (Getting worried...)
- At $\sigma_x/\sigma_{\text{range}} = 1.0$, the factor is $\lambda_{\beta} = 0.50$, corresponding to 50% BIAS!

At SNR of 10 or below, a regime common in astronomical observations, the bias reaches percent level or more. Several properties of this bias are worth emphasizing. First, the bias depends only on the ratio $\sigma_x/\sigma_{\text{range}}$, not on the absolute values of either quantity. Second, it is independent of sample size n and cannot be resolved by collecting more data. Even with infinite training samples, the bias remains. Third, the magnitude of attenuation bias depends only on errors in x , not errors in y . Errors in y add scatter to the predictions but do not systematically bias them, because positive and negative

errors average out over the dataset. In other words, obtaining better or more accurate labels does not solve the problem.⁴

This asymmetry between x and y arises because they play fundamentally different roles in regression: we predict y given x , not the reverse. Errors in x change the question being asked—we are predicting from a corrupted version of the input—which systematically weakens the apparent relationship. However, as we will see in Section 4, while σ_y does not affect the bias itself, it does play a role in correcting the bias: the Deming regression solution (Deming 1943; Golub & Van Loan 1950; Ting 2025c) and its neural network generalization require specifying both σ_x and σ_y to properly weight the loss terms.

Why this bias is excessively annoying:

- ☑ Getting MORE data doesn't fix it. (It's a bias, not variance!)
- ☑ It's caused by errors in the INPUTS (x), not the labels (y).



So, linear math is neat, but what about...

3. ATTENUATION BIAS IN NEURAL NETWORKS derived analytically. Paper I further argued that the same bias extends to polynomial and nonlinear models. However, for nonlinear models or high-dimensional inputs, the attenuation factor can no longer be expressed as a simple function of σ_x and σ_{range} . The relationship between measurement error and attenuation becomes more complex, depending on the specific functional form.

This complexity makes it difficult to predict and therefore analytically correct the bias in closed form for general machine learning models. Here we first demonstrate that multi-layer perceptrons (MLPs) exhibit the same qualitative behavior.

This complexity makes it difficult to predict and therefore analytically correct the bias in closed form for general machine learning models.

empirically that multi-layer perceptrons (MLPs) exhibit the same qualitative layer perceptrons (MLPs) exhibit the same qualitative behavior. MLPs, also known as fully connected neural networks (Gybenko 1988; Hornik et al. 1999), are the simplest neural network architecture: each neuron in one layer is connected to every neuron in the next layer, with nonlinear activation functions between layers. Despite their simplicity, MLPs remain widely used in astronomy (e.g., Odewahn et al. 1992; Storrie-Lambrad et al. 1992; Collister & Lahav

But the problem is still there! Our fancy deep learning models get lazy too.

Real-world examples of the damage:

Common in astro!

The math gets too hard here. No simple formula to save us.

OK, let's test this on Neural Networks! Do they get lazy too?

4

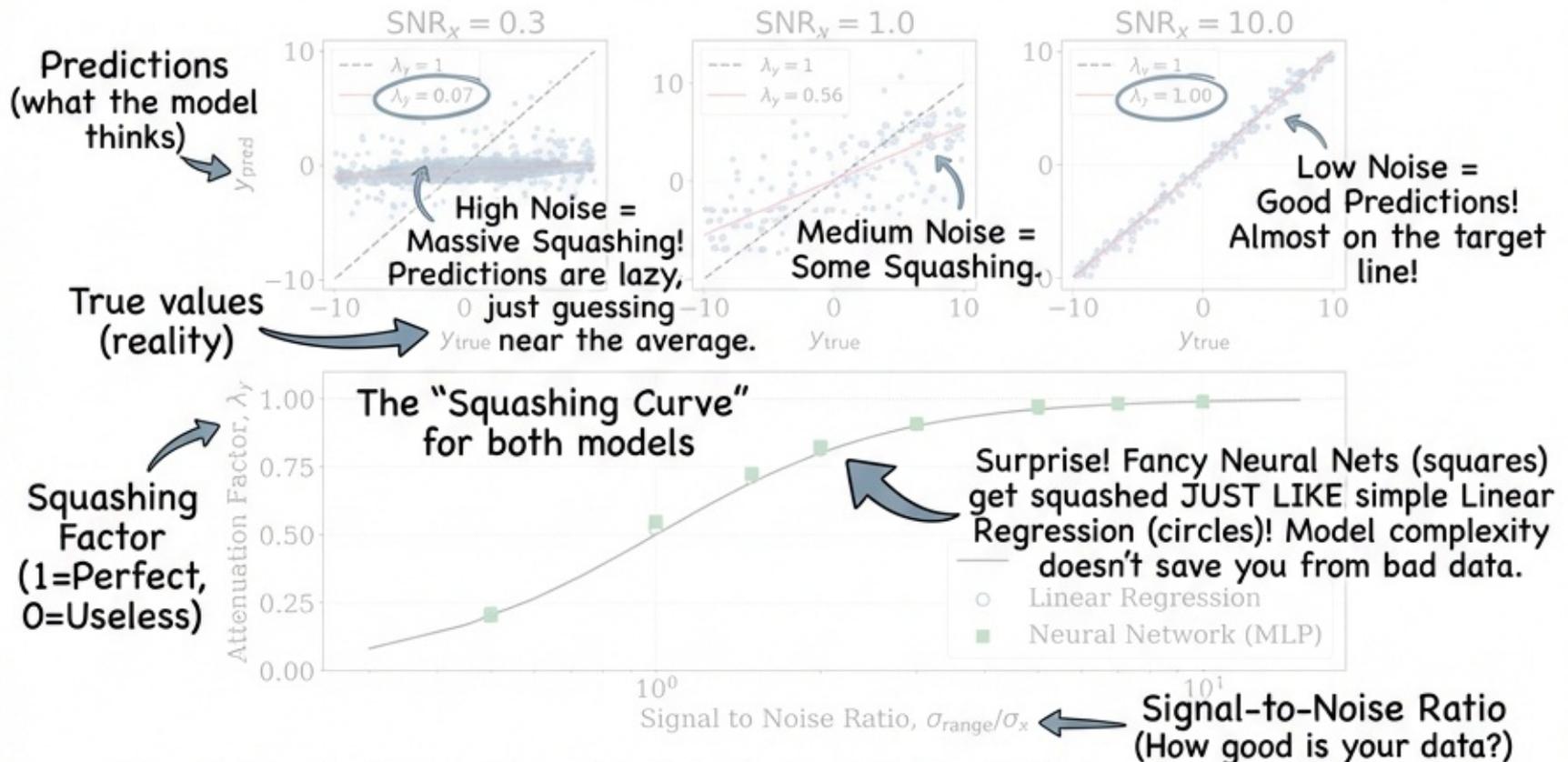


Figure 2. Attenuation bias in neural networks for the true function $y = 2x$. Top: Predicted versus true y on held-out test data for an MLP (2 hidden layers, 64 units) at three SNR_x values. At $SNR_x = 0.3$, predictions collapse to $\lambda_y \approx 0.1$; at $SNR_x = 1$, $\lambda_y \approx 0.5$; at $SNR_x = 3$, $\lambda_y \approx 0.9$. Bottom: Attenuation factor λ_y versus SNR_x for linear regression (blue open circles) and neural network (green filled squares). Both follow the theoretical curve $\lambda_y = 1/(1 + SNR_x^{-2})$ (black solid line), demonstrating that model complexity does not protect against attenuation bias.

2004; Guo et al. 2020; Wong et al. 2020; Wang et al. 2021; Li et al. 2022b; Vynatheya et al. 2022; Winecki & Kochanek 2024) and serve as the foundation for more complex architectures.

To test whether neural networks suffer from attenuation bias, we conduct numerical experiments following the setup of Paper I. We consider the true function $y = 2x$, chosen for direct comparison with the linear theory, and sample $N_{train} = 1000$ training values and $N_{test} = 200$ test values of x_{true} uniformly from $[-5, 5]$. We add Gaussian noise to obtain $x_{obs} = x_{true} + \delta_x$ where $\delta_x \sim \mathcal{N}(0, \sigma_x^2)$, and vary σ_x to test different SNR regimes. We define $SNR_x \equiv \sigma_{range}/\sigma_x$ as the signal-to-noise ratio in the input.

We train two models on each dataset. The first is ordinary least squares linear regression. The second is an MLP with 2 hidden layers of 64 units each and ReLU activations, trained for 20000 epochs using the Adam optimizer with learning rate 0.03.

For neural networks, we cannot directly measure a single coefficient β as we did for linear regression in Section 2. The MLP used here has over 4000 parameters distributed across multiple layers, rather than a single interpretable slope. Instead, we quantify attenuation

through the predictions themselves. We define λ_y as the slope of the regression of y_{pred} against y_{true} . If a model suffers from attenuation bias, its predictions will be compressed toward the mean: when we plot y_{pred} versus y_{true} , the slope will be less than unity. A value of $\lambda_y = 1$ indicates no attenuation, while $\lambda_y < 1$ indicates that predictions are systematically pulled toward the sample mean.

As discussed in Paper I, for linear regression in one dimension with true relationship $y = \beta x$, these two measures of attenuation are equivalent. The predicted value is $\hat{y} = \hat{\beta} x_{obs}$, and when we regress predictions against true values, the attenuation in the slope estimate $\hat{\beta}$ translates directly to attenuation in the predictions. In this case, $\lambda_y = \lambda_{\hat{\beta}} = 1/(1 + (\sigma_x/\sigma_{range})^2)$. For nonlinear models or multivariate inputs, λ_y remains a useful empirical measure of how much the model compresses the dynamic range of predictions.

Figure 2 shows the results. The top row shows MLP predictions on held-out test data at three specific SNR_x values: at $SNR_x = 0.3$ (severe noise), predictions are compressed to $\lambda_y \approx 0.1$; at $SNR_x = 1$ (noise comparable to signal range), $\lambda_y \approx 0.5$; and at $SNR_x = 3$, $\lambda_y \approx 0.9$. The bottom panel shows that both linear re-

Since NNs are complex black boxes, we measure the squashing by checking the slope of predictions vs reality.

Here's our experiment recipe: take a simple relationship, then corrupt the inputs with noise. Since NNs are complex black boxes, we measure the squashing by slope of predictions vs reality.

So, complex models don't save us. What now? Enter our proposed solution: LATENTNN!

gression and the MLP follow the theoretical attenuation curve $\lambda_y = 1/(1 + (\sigma_z/\sigma_{range})^2)$ across the full range of SNR_y. For the linear true function $y = 2x$, the MLP effectively learns a linear mapping, so it follows the same theoretical curve as linear regression. The important point is that our MLP here, despite having thousands of parameters and the flexibility to fit arbitrary nonlinear functions, exhibits the same systematic bias as the much simpler linear regression. As we argued in Paper I, model complexity does not protect against attenuation bias.

Regularization does not resolve the issue either. One might hope that L_2 regularization (weight decay), which penalizes large weights, could counteract the attenuation

The Core Idea: Don't treat noisy data as "truth". Instead, treat the *true, unknown* values as hidden ("latent") variables that we need to estimate along with the model!

(Tikhonov 1959; Carroll et al. 1995; Carroll et al. 2006; Bishop 2006). This systematic bias mismatch cannot correct. We verified this empirically by sweeping over a range of L_2 regularization strengths: the attenuation factor λ_y remains unchanged regardless of the regularization parameter.

4. LATENTNN: CORRECTING ATTENUATION BIAS WITH LATENT VARIABLES

Having established that neural networks suffer from attenuation bias, we now develop a method to correct it. The approach builds on a classical idea from the statistics literature on *errors-in-variables* regression (Carroll et al. 1995; Carroll et al. 2006; Bishop 2006): rather than treating noisy observations as exact, we treat the unknown true values as latent variables to be estimated alongside the model parameters. For linear regression, this leads to Deming regression (Deming 1943). Here we show that the same principle extends naturally to neural networks, visiding a practical correction method we call LatentNN.

4.1. The Latent Variable Solution in Linear Regression

Standard regression treats x_{obs} as exact, which leads to bias. An alternative approach recognizes that x_{true} is an unknown quantity for each data point. Rather than treating the observed values as truth, we can treat the true values as latent variables to be estimated alongside the model parameters.

Consider observation i . We assume that the observed input $x_{obs,i}$ is a noisy measurement of the true value $x_{true,i}$, and the observed output $y_{obs,i}$ is a noisy measurement of the true output $y_{true,i}$, which follows the

linear relationship $y_{true,i} = \beta x_{true,i}$. Mathematically,

$$x_{obs,i} = x_{true,i} + \delta_{x,i}, \quad \delta_{x,i} \sim \mathcal{N}(0, \sigma_x^2), \quad (3)$$

$$y_{obs,i} = \beta x_{true,i} + \delta_{y,i}, \quad \delta_{y,i} \sim \mathcal{N}(0, \sigma_y^2). \quad (4)$$

Given these assumptions, the joint probability of observing both $x_{obs,i}$ and $y_{obs,i}$ given the true value $x_{true,i}$ and the model parameter β is

Quick detour: You might think "regularization" (a common trick to prevent overfitting) would help. Nope! It fixes overfitting, not this noise bias problem.

Here x from a Gaussian distribution with mean μ and

Over all N observations, the full likelihood is

$\mathcal{L}(\beta, \{x_{true,i}\}) = \prod_{i=1}^N \mathcal{P}(x_{obs,i}, y_{obs,i} | x_{true,i}, \beta)$

The model tries to figure out the *true* point on the line that generated the messy data we actually see.



This is just writing down the problem mathematically: Observed = True + Noise.

This big equation is just asking: "Given the messy data we have, what's the most likely set of *true* values and the *true* line slope (β)?"

at the top level, the true values $x_{true,i}$ form the intermediate level, and the actual measurements ($x_{obs,i}, y_{obs,i}$) sit at the bottom level. The model constrains the intermediate level, while the intermediate level generates the observations.

This seems hard (too many unknowns!), but it works because all the "true" points must fall on a single straight line. That's a huge constraint that makes it solvable!

The linear case (Deming regression) has an exact, albeit scary-looking, mathematical solution...

defining the error variance ratio $r = \sigma_y^2/\sigma_x^2$, the estimator is

$$\hat{\beta}_{\text{Deming}} = \frac{S_{yy} - rS_{xx} + \sqrt{(S_{yy} - rS_{xx})^2 + 4rS_{xy}^2}}{2S_{yy}}, \quad (7)$$

where $S_{xx} = \sum_i (x_{\text{obs},i} - \bar{x})^2$, $S_{yy} = \sum_i (y_{\text{obs},i} - \bar{y})^2$, and $S_{xy} = \sum_i (x_{\text{obs},i} - \bar{x})(y_{\text{obs},i} - \bar{y})$ are the sum of squared deviations and cross-products, with \bar{x} and \bar{y} denoting the sample means. In the limit $\sigma_x \rightarrow 0$ (equivalently $r \rightarrow \infty$), this reduces to the ordinary least squares estimator $\hat{\beta}_{\text{OLS}} = S_{xy}/S_{xx}$. When $\sigma_x > 0$, the Deming estimator yields a steeper slope than OLS, correcting for the attenuation (see also Ting 2025b).

The solution for the optimal $x_{\text{true},i}$ provides further intuition. For fixed β , the optimal estimate is a weighted average

$$x_{\text{true},i} = w_x \cdot x_{\text{obs},i} + w_y \cdot \frac{y_{\text{obs},i}}{\beta}, \quad (8)$$

where the weights are (see Appendix B for derivation)

$$w_x = \frac{r}{r + \beta^2}, \quad w_y = \frac{\beta^2}{r + \beta^2}. \quad (9)$$

The weights sum to unity. When σ_x is small compared to σ_y/β , the intuition is beautiful: it's a smart **WEIGHTED AVERAGE**. Trust the measurement that is **less noisy**.

NOW, THE KEY INSIGHT: The "linear" part wasn't essential. The hierarchical structure was.

4.2. The LatentNN Formulation

The linear form of Deming regression was not essential—what mattered was the hierarchical structure where latent values generate observations and the model constrains the latent values. The linear form permits an analytic solution, but the underlying principle of joint optimization over model parameters and latent true values extends directly to neural networks.

We replace the linear model βx with a neural network $f_\theta(x)$, where θ represents all the network weights and biases. The hierarchical structure remains the same: the network parameters θ sit at the top level (analogous to β), the latent values $x_{\text{latent},i}$ form the intermediate level (analogous to $x_{\text{true},i}$), and the observations $(x_{\text{obs}}, y_{\text{obs}})$ sit at the bottom level. The likelihood has exactly the

same form as in the linear case,

$$\mathcal{L}(\theta, \{x_{\text{latent},i}\}) = \prod_{i=1}^N \mathcal{N}(x_{\text{obs},i} | x_{\text{latent},i}, \sigma_{x,i}^2) \times \mathcal{N}(y_{\text{obs},i} | f_\theta(x_{\text{latent},i}), \sigma_{y,i}^2), \quad (10)$$

with $f_\theta(x_{\text{latent},i})$ replacing $\beta x_{\text{true},i}$. Here $\sigma_{x,i}$ and $\sigma_{y,i}$ denote the per-sample uncertainties across observations.

Taking the negative log-likelihood and absorbing constants, we get

$$\mathcal{L}(\theta, \{x_{\text{latent},i}\}) = \sum_{i=1}^N \frac{(y_{\text{obs},i} - f_\theta(x_{\text{latent},i}))^2}{\sigma_{y,i}^2} + \sum_{i=1}^N \frac{(x_{\text{obs},i} - x_{\text{latent},i})^2}{\sigma_{x,i}^2}. \quad (11)$$

Comparing to the Deming regression objective in Appendix B, we see the same structure: a prediction loss term weighted by $1/\sigma_{y,i}^2$ and a latent likelihood term

This leads to a new "loss function" to minimize, with the same two parts:

- (1) Match observations
- (2) Keep latents close to data

but infer the true value from the model when the measurement is noisy. The uncertainties $\sigma_{x,i}$ and $\sigma_{y,i}$ must be specified for proper correction.

The optimization problem is to find

$$\hat{\theta}, \{\hat{x}_{\text{latent},i}\} = \arg \min_{\theta, \{x_{\text{latent},i}\}} \mathcal{L}(\theta, \{x_{\text{latent},i}\}). \quad (12)$$

For a dataset with N samples and p input dimensions, this involves optimizing over the network parameters θ plus $N \times p$ latent variables. The total number of parameters is thus much larger than in standard neural network training, scaling with the training set size.

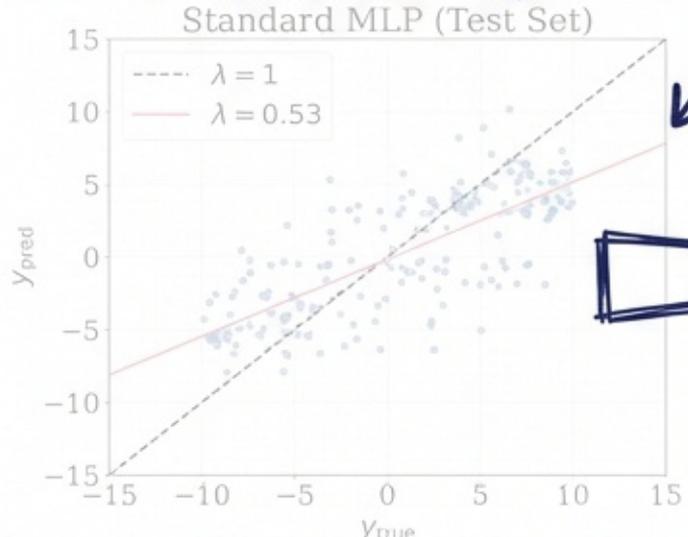
The key difference from the linear case is that Equation 12 has no closed-form solution. In Deming regression, the linear structure allows us to solve for both $\hat{\beta}$ and $\{\hat{x}_{\text{true},i}\}$ analytically (see Appendix B). For a neural network f_θ , the nonlinear dependence on θ precludes an analytic solution, and we must resort to numerical optimization.

4.3. Implementation

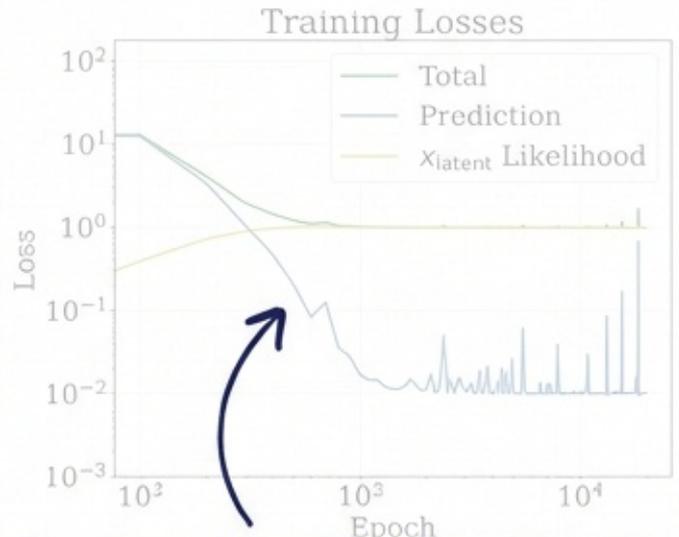
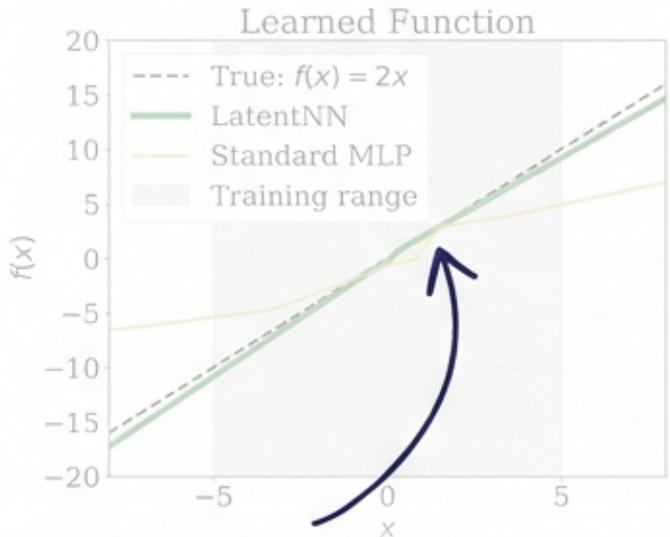
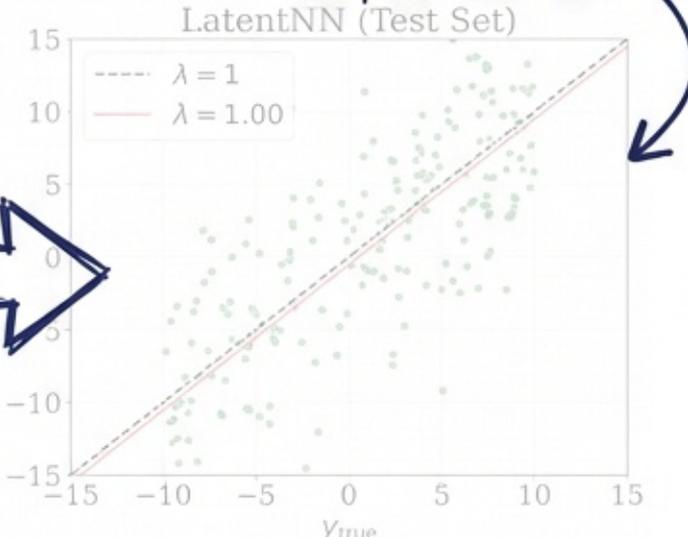
We solve Equation 12 using gradient descent, the standard optimization algorithm for neural networks. Let

THE CATCH: Unlike the linear case, there's no simple formula for the answer. We have to use numerical optimization (like gradient descent) to find it.

6 See? The standard Neural Net $CI(MLP)$ is still SQUASHED by the noise! Slope is wrong.



But our LatentNN fixes it! The predictions line up with the truth! Slope is correct.



This plot shows the actual functions they learned. The true function is $f(x)=2x$ (dashed). The standard one is too flat, but LatentNN nailed it!

This shows the training in action. The model is balancing two goals: fitting the data (Prediction loss \downarrow) and keeping latent values plausible (x_{latent} Likelihood \uparrow). It's a tug-of-war!

$\Phi = (\theta, \{x_{latent,i}\})$ denote the full set of parameters to be optimized. At each iteration t , we update

$$\Phi^{(t+1)} = \Phi^{(t)} - \eta \nabla_{\Phi} \mathcal{L}, \quad (13)$$

where η is the learning rate and $\nabla_{\Phi} \mathcal{L}$ is the gradient of the loss with respect to all parameters. This is the same update rule used in standard neural network training, but with Φ now including both the network weights θ and the latent inputs $\{x_{latent,i}\}$. This represents a conceptual shift from typical neural network applications, where we optimize only over θ . Here, we treat the observed inputs as noisy measurements rather than ground truth, and optimize jointly over θ and the latent inputs.

In modern deep learning frameworks such as PyTorch or JAX, this is straightforward to implement.

Here's that "numerical optimization" we mentioned! It's just standard gradient descent, but we update everything at once...

...by treating the latent values as just more parameters to learn! Clever trick.

Here's the SYSTEMATIC PROOF.
LatentNN vs. Standard MLP across ALL noise levels.

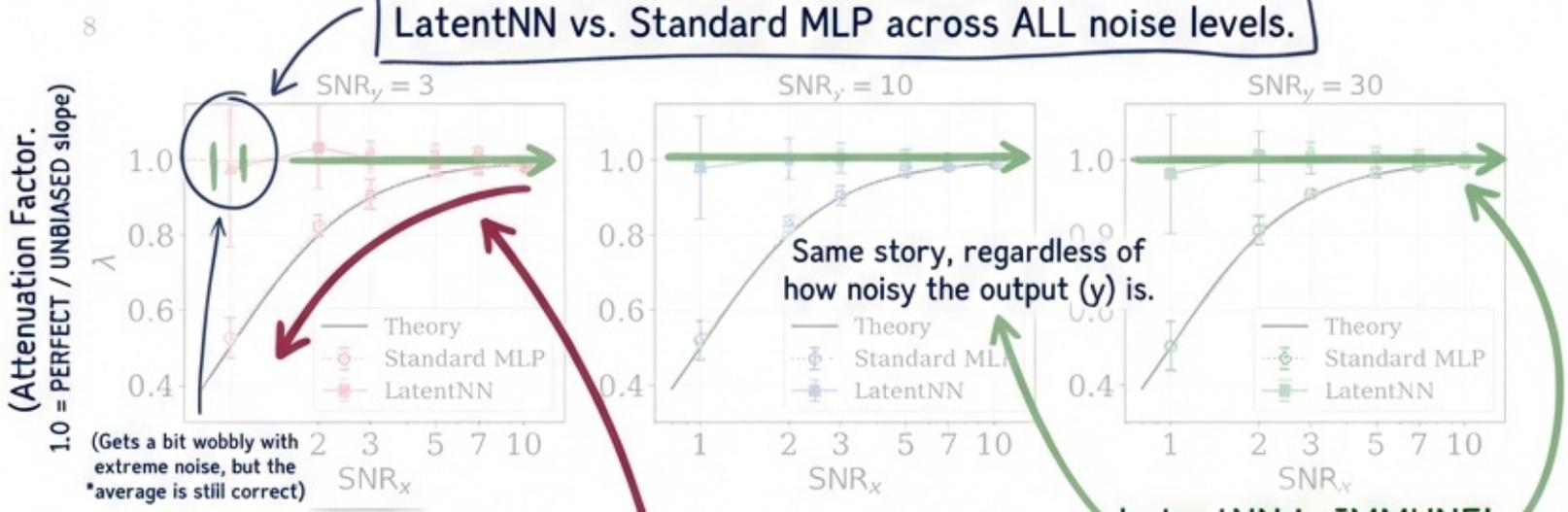


Figure 4. (Input Signal-to-Noise. SNR_x for three SNR_y levels (3, 10, 30; different deviation. Left = VERY NOISY inputs) seeds, using a three-fold data split (train, validation, test) and hyperparameter selection. Standard MLP (open circles, dotted) follows the theoretical prediction of SNR_y . LatentNN (filled squares, solid) maintains mean $\lambda_y \approx 1$ across all tested SNR_x . LatentNN achieves unbiased predictions on average, though with increased run-to-run scatter at lower SNR_x as the problem becomes increasingly ill-conditioned.

The Standard Net gets "squashed" by noise, exactly as theory predicts. It fails hard here! →

of network parameters, the configurations where the network overfits. We compare different configurations rather than learning. Weight decay on the network parameters prevents this by penalizing large weights, encouraging smoother functions that generalize better. Weight decay is applied only to θ , not to the latent values x_{latent} —the latent values are already constrained by the latent likelihood term.

shows a representative comparison at $SNR_x = 1$, using $SNR_y = 10$ as in Section 3. At $SNR_x = 1$, the measurement error equals the signal range ($\sigma_x = \sigma_{range}$)—an extreme regime where the noise is 100% of the data's dynamic range. The top panels show predicted versus true values on the held-out test set. The standard MLP shows $\lambda_y \approx 0.5$, demonstrating significant bias. LatentNN recovers $\lambda_y \approx 1$, showing unbiased predictions.

5. NUMERICAL EXPERIMENTS

Time for a rigorous, fair test...

We now validate the LatentNN approach on test cases where the true relationship and measurement errors are known, allowing direct comparison between estimated and true attenuation factors. We begin with the one-dimensional case then extend to multivariate inputs with correlated features.

3.1. One-Dimensional Case

We use the same data setup and network architecture as in Section 3: the true relationship $y = 2x_{true}$ with 1000 training samples, and an MLP with 2 billion parameters of 64 units trained for 20000 epochs. To prevent overfitting, we add a learnable parameter λ to the output layer, initialized to $x_{obs,i}$.



To ensure fair comparison with no data leakage, we use a three-fold data split: 1000 samples for training, 200 for validation (hyperparameter tuning), and 200 for final evaluation. The validation set is used exclusively for selecting hyperparameters—for LatentNN we maximize λ_y on the validation set, while for the standard MLP we minimize validation loss. Final performance metrics are computed on the held-out test set, which is never

the bottom-left panel shows the learned functions: the standard MLP learns an attenuated slope, while LatentNN recovers the true relationship $f(x) = 2x$.

The bottom-right panel shows the training dynamics for LatentNN. At initialization, $x_{latent} = x_{obs}$, so the x_{latent} likelihood term starts near zero. As training proceeds, the prediction loss decreases as the network learns the mapping, but the x_{latent} likelihood increases. This increase is expected and desirable: to achieve accurate predictions, the latent values must shift from the noisy values x_{obs} to values x_{true} , which are closer to the true relationship. The total loss decreases as the network learns the true relationship, resulting in improved prediction accuracy on the test set. At convergence, the optimization has found the balance point where the latent values have moved sufficiently toward their true positions to enable unbiased prediction.

We use a strict data split to prevent "cheating" (data leakage).

To systematically evaluate performance, we repeat each experiment 8 times with different random seeds, using the same three-fold data split (train/validation/test).

CRITICAL DETAIL: They are tuned differently! Standard MLP just optimizes for low error, while LatentNN is specifically tuned to get the slope (λ) correct on validation data.

The secret sauce: The two loss terms balance each other differently depending on the noise level.

ing three SNR_y values (3, 10, and 30). As expected from the theory in Section 2, varying SNR_y has little effect on attenuation—the bias depends primarily on input noise. The standard MLP follows the theoretical attenuation curve $\lambda_y = 1/(1 + (\sigma_x/\sigma_{\text{range}})^2)$ regardless of SNR_y . LatentNN maintains $\lambda_y \approx 1$ across all tested SNR_x and SNR_y values. LatentNN provides robust correction even at $\text{SNR}_x = 1$, though with increased run-to-run scatter at lower SNR_x values as the problem becomes increasingly ill-conditioned.

The mechanism of correction differs between SNR regimes. At high SNR where $\sigma_x/\sigma_{\text{range}} \ll 1$, the x_{latent} likelihood term in Equation 11 dominates and the latent values remain close to the observations ($x_{\text{latent}} \approx x_{\text{obs}}$). At moderate SNR ($\text{SNR}_y \approx 2-3$), the two loss terms compete more evenly, and weight decay tuning becomes important for LatentNN. The standard MLP, by contrast, fails regardless of weight decay—it lacks the mechanism to correct attenuation bias, and adjusting regularization only trades off between overfitting and underfitting without addressing the systematic bias.

LatentNN generalizes well to held-out test data. Weight decay regularization applied only to the network parameters θ prevents the network from overfitting to the specific denoised distribution. The test set performance is consistent across all SNR_x values, indicating that the learned function f_θ captures the true underlying relationship rather than artifacts of the training data.

Now for the REAL WORLD stuff: Multivariate Inputs!

Term 1: Minimize prediction error on output (y).

Term 2: Constrain latent X to be plausible given noisy input.

The full multivariate loss function.

5.2. Multivariate Inputs

Real astronomical applications involve multivariate inputs. The LatentNN formulation generalizes to p input dimensions. The latent variables $X_{\text{latent},ij}$ form an $N \times p$ matrix, where p is the number of input features. The loss function from Equation 11 extends naturally to

$$\mathcal{L} = \sum_{i=1}^N \frac{(y_i - f_\theta(X_{\text{latent},i}))^2}{\sigma_{y,i}^2} + \sum_{i=1}^N \sum_{j=1}^p \frac{(X_{\text{obs},ij} - X_{\text{latent},ij})^2}{\sigma_{x,ij}^2}, \quad (14)$$

where $\sigma_{y,i}$ is the output uncertainty for sample i and $\sigma_{x,ij}$ is the input uncertainty for sample i and dimension j . For simplicity, we assume homogeneous noise σ_x and σ_y across all samples and dimensions, which allows us to visualize results as a function of a single SNR_x .

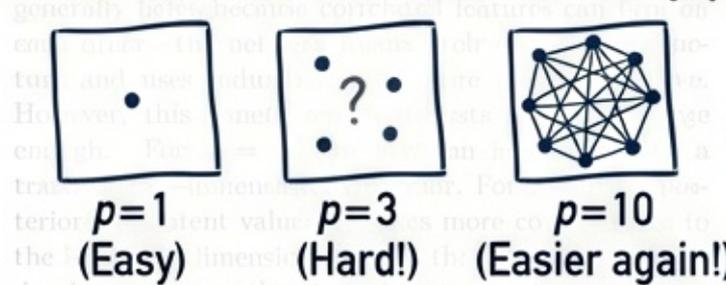
As shown in Paper I, for uncorrelated features each dimension attenuates independently, recovering the one-dimensional result. We therefore focus on correlated features, which better represent astronomical data like

stellar spectra where different pixels respond coherently to changes in stellar parameters. We generate $N_{\text{train}} = 1000$ training samples where all features derive from a single latent variable: $x_j = a_j z$ where $z \sim \text{Unif}(-0.5, 0.5)$ and a_j varies linearly from 0.2 to 0.4 across dimensions. The true output is $y = x^T a$, centered at zero. We fix $\text{SNR}_y = 10$ as in the single-variate experiments. We use the same network architecture and training settings as Section 3, with the same three-fold data split (1000 train, 200 validation, 200 test) and weight decay grid search strategy. We test $p = 3, 10, 30$ dimensions at SNR_x values from 1 to 10, repeating each experiment 8 times.

Figure 5 shows the results. Paper I derived an analytic expression for the attenuation factor in correlated linear regression: $\lambda_y = \lambda_\beta = \sum_j a_j^2 / (1/\text{SNR}_x^2 + \sum_j a_j^2)$. The gray curves show this theoretical prediction. Just as in the single-variate case, the standard MLP closely follows this linear regression theory curve, confirming that additional model complexity does not mitigate the fundamental attenuation bias. With correlated features, attenuation is weaker than the $p = 1$ case because independent noise in different dimensions must conspire coherently to displace observations along the signal direction. This bias is mitigated compared to smaller p , but remains significant at lower SNR.

The results show a fascinating “Goldilocks” effect with dimensionality (p):

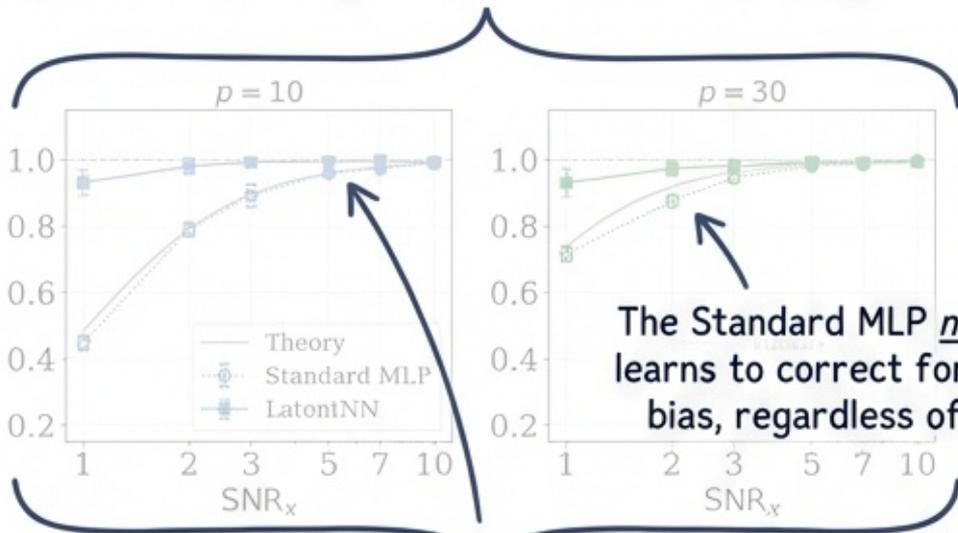
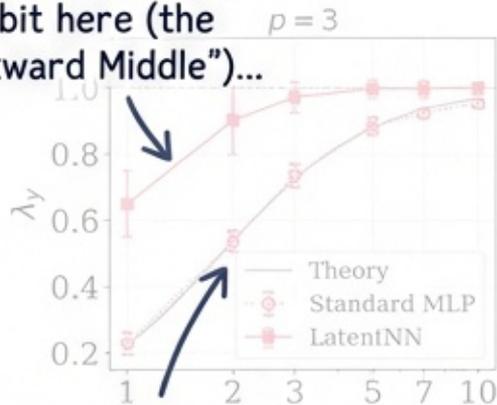
At high p , correlated features “lean on each other”, providing redundancy that helps constrain the latent values. It’s actually easier than intermediate p !



At high p , correlated features “lean on each other”, providing redundancy that helps constrain the latent values. It’s actually easier than intermediate p !

Here's the visual proof of the "tug-of-war" from the last page!

LatentNN struggles a bit here (the "Awkward Middle")...



...but still beats the Standard MLP, which is stuck on the biased theory curve.

Strength in Numbers! With more features ($p=10, 30$), LatentNN's correction mechanism works almost perfectly ($\lambda_y \approx 1$).

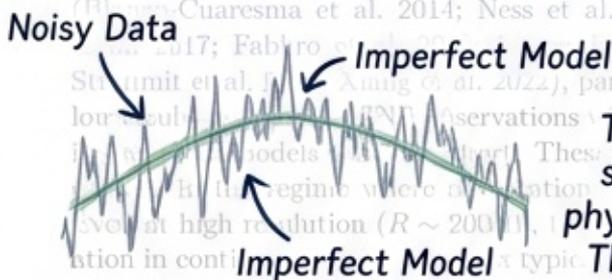
challenging $p = 3$ regime, LatentNN substantially outperforms the standard MLP: at $\text{SNR}_x = 1$, LatentNN achieves $\lambda_y \approx 0.7$ compared to the standard MLP's $\lambda_y \approx 0.25$, more than tripling the recovery range.

Real applications of LatentNN to survey data are underway, but this paper focuses on introducing the technique and establishing its theoretical foundation as an extension of Deming regression to neural networks. Here we focus on inferring overall metallicity $[M/H]$ for clusters such as individual elemental abundances.

OK, enough toy models. Let's try a REAL astronomical problem!

Stellar spectra represent the primary astronomical application motivating this work. Data-driven models have been applied extensively to spectroscopic surveys (Blanton & Cuadras 2014; Ness et al. 2015; Ting et al. 2017; Fabro et al. 2019; Bovy 2019; Strumit et al. 2024), particularly for low-resolution observations. These models are trained in the regime where attenuation bias is not a concern, typically at high resolution ($R \sim 2000$), in contrast to typical spectroscopic applications with $\sigma_{\text{range}} \approx 5\text{--}20\%$ ($\sigma_{\text{range}} \approx 0.05\text{--}0.2$).

For the weaker lines of $\sigma_{\text{range}} \approx 0.05$, a spectrum with conventional $\text{SNR} = 20$ (defined relative to the normalized flux) has per-pixel uncertainty $\sigma_x = 0.05$, giving $\text{SNR}_s = \sigma_{\text{range}}/\sigma_x \approx 1$ in our definition; conventional $\text{SNR} = 100$ corresponds to $\text{SNR}_s \approx 10$. We emphasize that what matters for attenuation bias is this ratio of signal range to measurement uncertainty, not the conventional spectroscopic SNR. Since the ratio $\sigma_{\text{range}}/\sigma_x$ is often of order 10 or less for many spectral pixels in typical spectroscopic applications, they fall squarely in the regime where attenuation bias becomes important—especially for elements with weak or limited spectral features.



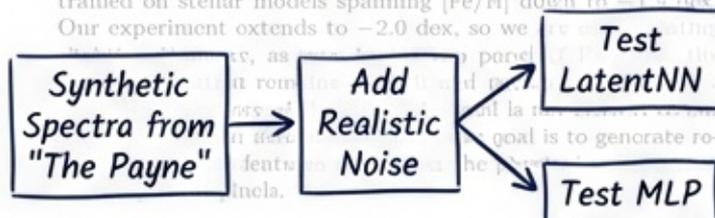
The challenge: Real stellar spectra are noisy, and our physical models aren't perfect. This combination is a recipe for attenuation bias!

We use the pretrained synthetic spectral emulator with The Payne (Ting et al. 2019) in The Payne github repository⁵ to generate spectra at resolution $R \approx 22500$, following Paper I. The underlying spectral model uses atomic and molecular line data from Kurucz (1970, 1983, 2005, 2017).⁶ We focus on a metallicity-sensitive region around 16200 \AA .

Now the number of informative features is small, and to test correction and correction, we test three pixel choices represent the range of spectral choices for different labels. While major elements like Fe have many strong lines, attenuation bias is most problematic for elements with only a few spectral features (e.g., K, V; see fig. 7 and 8 of Xiang et al. 2019).

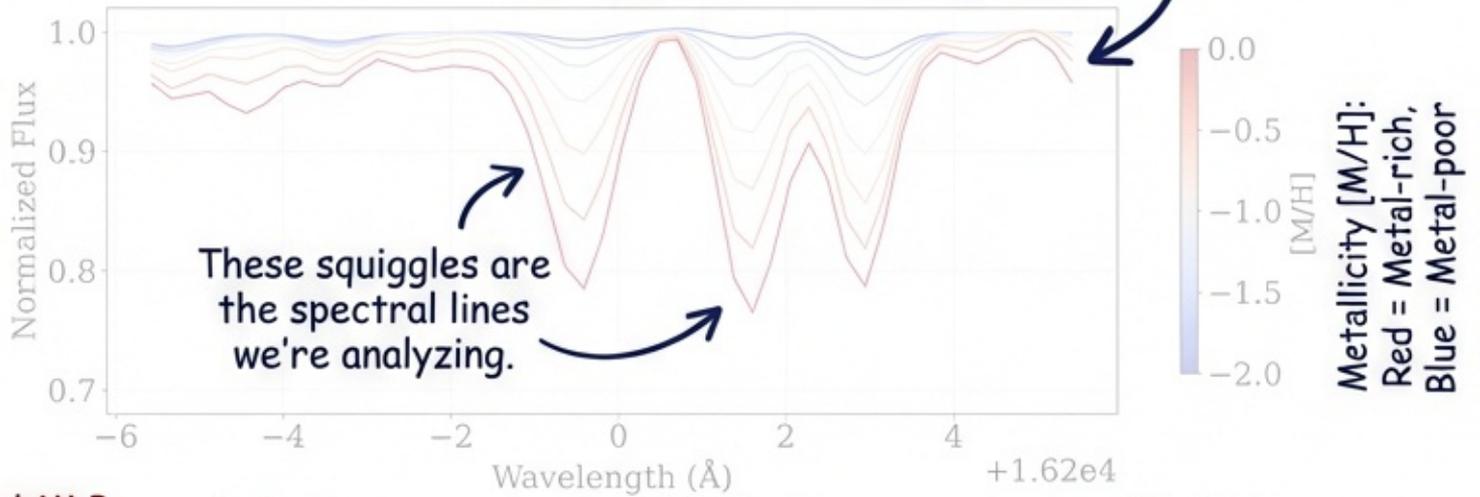
⁵ <https://github.com/tingyuansen/The.Payne>

⁶ This emulator from the github repository of The Payne is trained on stellar models spanning $[\text{Fe}/H]$ down to -1.5 dex. Our experiment extends to -2.0 dex, so we use the -1.5 dex model.



We're using a sophisticated emulator to create realistic test data, focusing on tricky elements like Iron (Fe).

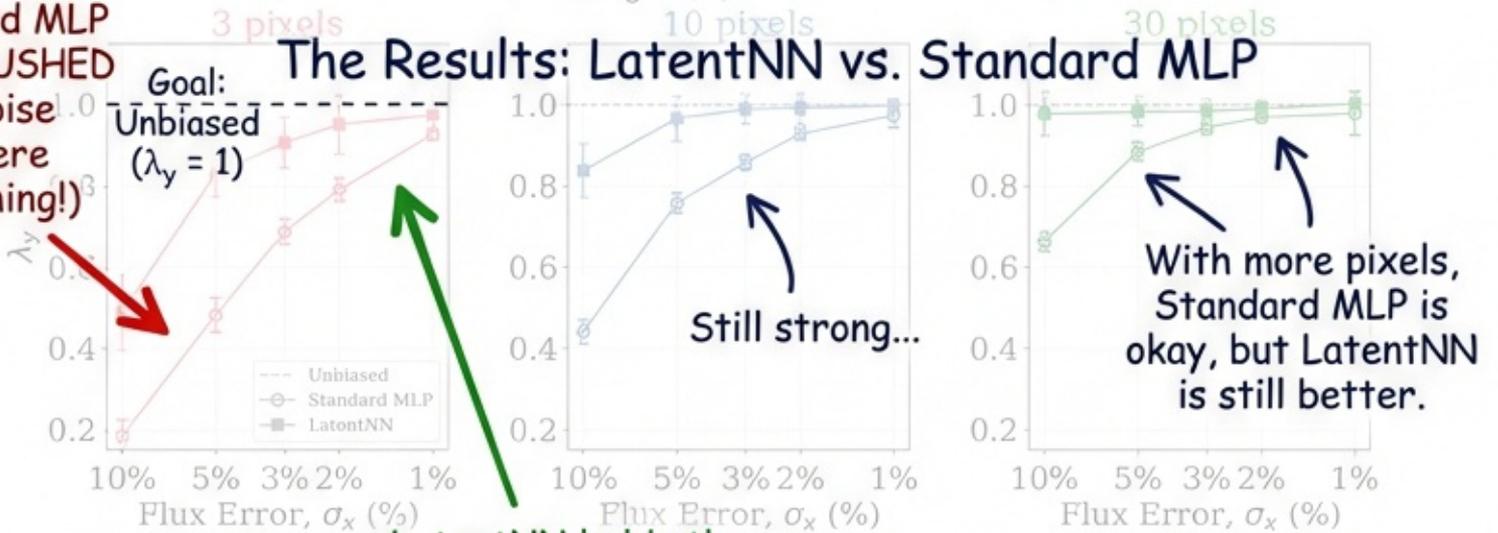
THE FINAL SHOWDOWN: Real(ish) Astronomy Data!



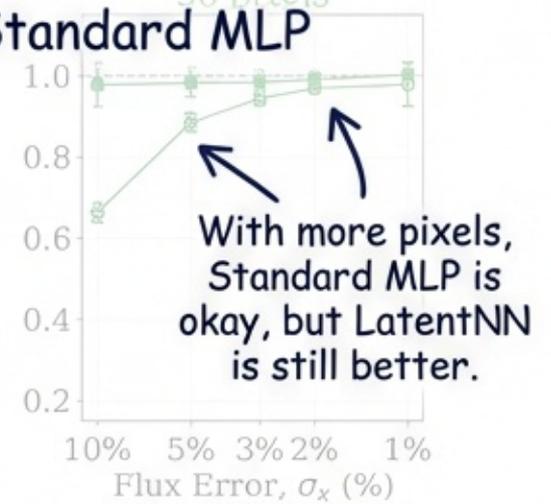
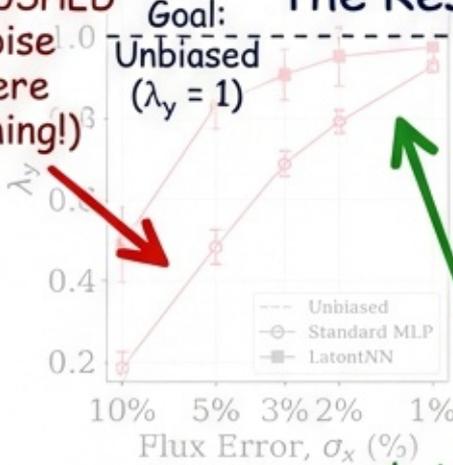
These squiggles are the spectral lines we're analyzing.

Metallicity [M/H]:
Red = Metal-rich,
Blue = Metal-poor

Standard MLP gets CRUSHED by noise (severe squashing!)



The Results: LatentNN vs. Standard MLP



LatentNN holds the line! Near-perfect recovery ($\lambda_y \approx 1$).

This is why the '3 pixels' plot is so important! Real data is often sparse.

Figure 6. Application to stellar spectra. Top: Continuum-normalized flux variation with metallicity [M/H]. Bottom: Attenuation factor λ_y vs. flux error σ_x for 3, 10, and 30 pixels. Error bars show standard deviation over 8 runs. The standard MLP shows significant attenuation that worsens with fewer pixels and higher noise. LatentNN (filled squares) maintains $\lambda_y \approx 1$ across all noise levels. Even for 3 pixels where correction is more challenging, LatentNN also maintains the attenuation bias compared to the standard MLP (e.g., $\lambda_y \approx 0.5$ vs. 0.2 at 10% flux error).

In practice, even surveys with thousands of pixels have only a small number of pixels that carry information for a given label. Adding realistic noise to challenge the models.

We generate samples spanning [M/H] from -2.0 to 0.0 dex, with fixed stellar parameters ($T_{\text{eff}} = 4750$ K, $\log g = 2.5$, typical of red clump stars). All abundances scale with [M/H] for simplicity. We use the same three-fold data split as in the synthetic experiments (1000 train, 200 validation, 200 test) with hyperparameters selected on the validation set and final evaluation on the held-out test set. We add Gaussian noise with standard deviation σ_s to each flux pixel and $\sigma_{[M/H]} = 0.05$ dex to the metallicity labels, representing conservative uncertainties in training labels from surveys like APOGEE (Holtzman et al. 2015; García Pérez et al. 2016). We test flux error values $\sigma_s \in \{1\%, 2\%, 3\%, 5\%, 10\%\}$ of the normalized flux. We use the same network architecture

Figure 6 shows the results. The top panel illustrates how the spectral flux varies with [M/H] across the selected wavelength region, with color indicating metallicity from -2.0 (blue) to 0.0 dex (red). Feature depths range from 5–20%, but even features with 10–20% total depth do not give $\text{SNR}_x = 2-4$ at 5% flux error as one might naively expect. This is because spectral features vary approximately quadratically with metallicity (Rix et al. 2016). In Paper I, we showed that the effective SNR is

Conclusion: LatentNN reliably corrects attenuation bias, even with noisy, limited data typical of real surveys. It works!

The bottom panels show λ_y versus flux error for each pixel configuration. The behavior mirrors the synthetic multivariate experiments (Section 3.2): the standard MLP exhibits attenuation that worsens with fewer pixels and higher noise, while LatentNN provides robust correction when sufficient pixels are available. For 10 and 30 pixels, LatentNN maintains $\lambda_y \gtrsim 0.95$ across all tested noise levels. For 3 pixels, however, correction degrades at high noise: $\lambda_y \approx 0.5$ at 10% flux error and $\lambda_y \approx 0.8$ at 5% flux error. This matches the intermediate-dimensionality challenge discussed in Section 5.2—with few pixels, the posterior over latent values is complex yet there is insufficient redundancy to leverage. Even so, LatentNN outperforms the standard MLP across all configurations: at 10% flux error with 3 pixels, LatentNN achieves $\lambda_y \approx 0.5$ compared to the standard MLP's $\lambda_y \approx 0.2$.

For spectral features with $\sigma_{\text{range}} \approx 0.10$ -0.20 (strong lines), flux errors of $\sigma_x \lesssim 3$ -5% yield effective $\text{SNR}_x \gtrsim 2$ after accounting for the factor of ~ 2 reduction from quadratic features—consistent with our findings throughout that LatentNN works well for $\text{SNR}_x \geq 2$. For weaker features with $\sigma_{\text{range}} \approx 0.03$, correspondingly lower flux errors ($\sigma_x \lesssim 1$ -2%) are needed. At higher noise levels with few informative pixels, careful hyperparameter tuning becomes essential, and users should expect some residual bias.

7. DISCUSSION

These huge surveys operate 'exactly' in the high-noise regime where this bias is devastating. They are all affected.

Even at APOGEE resolution, attenuation bias is non-negligible at moderate SNR. For lower-resolution surveys where individual spectral features are blended and per-pixel σ_{range} is smaller, the ratio $\sigma_x/\sigma_{\text{range}}$ is larger and attenuation bias is more severe. Crucially, σ_{range} scales approximately proportionally with spectral resolution R : a survey like LAMOST ($R \sim 1800$) has σ_{range} roughly 10 times smaller than APOGEE ($R \sim 22800$). This means that even at high conventional $\text{SNR} \sim 100$ ($\sigma_x \sim 1\%$), the effective SNR_x remains in the regime where attenuation bias is significant.

THE REAL-WORLD CONSEQUENCE: We misclassify stars!



True:
Metal-Poor
([M/H] = -2)



Biased Prediction:
Metal-Rich
([M/H] \approx -1.5)

This bias makes rare, ancient stars look normal, totally messing up our understanding of the Milky Way's history!

systematically predicted to be more metal-rich, making the rarest objects harder to identify. This bias directly affects the inferred metallicity distribution function (MDF), particularly its metal-poor tail, which is crucial for constraining the formation history of the proto-Milky Way (Rix et al. 2022; Chen et al. 2024, 2025). If attenuation bias compresses the metallicity scale, the true slope of the metal-poor MDF tail would be steeper than observed, with implications for inferred gas masses, star formation efficiencies, and inflow histories in galactic chemical evolution models. Stellar ages suffer from compression at both ends, with the oldest stars underestimated and the youngest overestimated. Distances are similarly affected, leading to systematic errors in the inferred geometry and dynamics of the Milky Way. Obtaining better training labels does not solve the problem because the bias arises from noisy inputs, not labels.

When information is concentrated in sparse spectral features such as individual absorption lines, attenuation is most severe. Consider an element like K or V with only 2-3 usable lines, compared to Fe with tens of strong features: K and V abundances are effectively in the low-dimensional regime where σ_{range} is large, while Fe benefits from the high-dimensional compressed regime. This pattern—elements with fewer lines suffering larger systematic errors at extreme values—has been observed in practice (Xiang et al. 2019; Zhang et al. 2024) but is often attributed to data quality or model overfitting. LatentNN offers a way out, providing a path to unbiased science from noisy spectral data. Case closed!

**THE FINAL VERDICT:
It's not bad data or bad models ...
it's Attenuation Bias.**

The errors-in-variables problem has a long history in statistics (Deming 1919; Fuller 1987). Deming regression offers a way out, providing a path to unbiased science from noisy spectral data. Case closed!



Applications for Spectroscopic Surveys

approach of transferring abundances from high-resolution resolution surveys to LAMOST and Gaia XP (Andrae et al. 2020; Fallows & Sanders 2024; Khalatyan et al. 2024; Hattori 2025; Yang et al. 2025) has become standard practice—but these applications operate precisely in the regime where attenuation bias is severe.

Putting it in context:
How does LatentNN
compare to other ideas?

... (Kubie & Van Loan 1980; ... Vandewalle 1991) treat the true input values as latent variables to be estimated alongside the regression coefficients—but this approach has seen limited application in astronomy (but see Kelly 2007), in part because astronomical relationships are often non-linear and standard neural network frameworks do not handle them well. LatentNN extends the classical latent variable approach to neural networks, providing a principled framework for astronomical applications where input uncertainties are non-negligible.

Our LatentNN formulation can be understood as the maximum a posteriori (MAP) limit of a hierarchical Bayesian model. In the full Bayesian formulation, the true values x_{true} are latent variables with prior distributions centered on the observed values x_{obs} , and we would infer the posterior distribution over both model parameters θ and all the latent values. The LatentNN optimization finds the mode of this posterior. Neural networks provide a flexible function approximator that avoids the need for explicit feature engineering, while the latent variable framework provides the statistical foundation for handling

This formulation stands as a paradigm shift in understanding why correlated features exist, as discussed in Section 3.2. The correlated structure of the inputs, in turn, provides a natural structure to constrain the latent values: each latent value is informed not only by its own observed value, but also by what the learned model expects for that feature.

It is important to note that the latent values x_{latent} should not be interpreted strictly as “denoised” inputs. In high-dimensional settings with large p , the posterior distribution over x_{latent} can be broad and complex—many combinations of latent values may be consistent with the observed data and the learned model. What matters for our purposes is that the network parameters

The principle is flexible! Here's how it works with more complex noise (heteroscedastic).

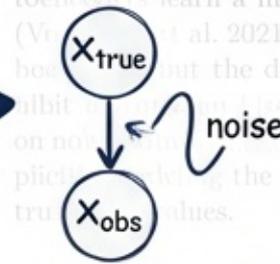
GENERALIZATIONS

And for classification tasks! Just swap out the loss function (e.g., CrossEntropy) but keep the latent variable part.

$p = 3$ regime in Figure 5), case-by-case hyperparameter tuning may be needed, but a systematic exploration of such regimes is beyond the scope of this paper.

This approach differs from common alternative strategies in astronomical research. Training with noise-augmented data (Eshbro et al. 2018) improves robustness to outliers but does not address bias in the mean prediction. Denoising autoencoders learn a mapping from noisy to clean inputs (Vahdat et al. 2021; Ghahramani & Van den Broek 2000; Doherty et al. 2019) but do not address the relationship between the observed and true values.

These common tricks help with robustness but DON'T fix the fundamental BIAS. We address the root cause.



The core idea: Treat true values as “hidden” variables we need to infer.

7.3. Generalizations

The core idea of maximising a likelihood with latent variables extends to many settings beyond the homoscedastic Gaussian case demonstrated here. As shown in Equation 14, heteroscedastic uncertainties are modeled by varying the weights of the likelihood term.

AHA MOMENT!

LatentNN isn't just “denoising” inputs. It's using a principled Bayesian framework to find the most likely true values given the noisy observations and the learned model.

The network parameters learn the correct relationship, even if individual latent values aren't perfect. That's the key!

For spectra due to emission lines, the latent variables are the true fluxes, and the observed values are the measured fluxes with heteroscedastic uncertainties. The likelihood term is then

$$\mathcal{L} = -\log p(\mathbf{x}_{\text{obs}} | \mathbf{x}_{\text{latent}}, \theta) = -\log p(\mathbf{y}_{\text{obs}} | f_{\theta}(\mathbf{x}_{\text{latent}})) \quad (15)$$

For Poisson noise, outliers, or other non-Gaussian errors, the Gaussian likelihood can be replaced with the appropriate distribution, giving $\mathcal{L} = -\log p(\mathbf{x}_{\text{obs}} | \mathbf{x}_{\text{latent}}) - \log p(\mathbf{y}_{\text{obs}} | f_{\theta}(\mathbf{x}_{\text{latent}}))$ for whatever form p takes.

For classification tasks with noisy inputs, the same latent variable approach applies. The prediction term becomes cross-entropy loss instead of mean squared error, while the latent likelihood term remains unchanged:

$$\mathcal{L} = \text{CrossEntropy}(f_{\theta}(\mathbf{x}_{\text{latent}}), \mathcal{Y}) + \sum_{i=1}^N (\mathbf{x}_{\text{obs},i} - \mathbf{x}_{\text{latent},i})^T \Sigma_{\mathbf{x},i}^{-1} (\mathbf{x}_{\text{obs},i} - \mathbf{x}_{\text{latent},i}) \quad (16)$$

Extension to other neural network architectures requires only that the network be differentiable with respect to the latent variables.

Bottom line: LatentNN gives us a rigorous, general way to handle noisy data and get unbiased science out of our surveys. Let's put it to work!

its inputs. Convolutional neural networks for imaging data, transformers for sequential data, or any other architecture can be used by simply treating the input as learnable parameters with an appropriate latent likelihood term.

7.A. Limitations

The method requires specifying both input uncertainties σ_x (or $\sigma_{x,i}$ in the heteroscedastic case) and output uncertainties σ_y . In practice, σ_y may need to be estimated based on the noise model of the instrument, while σ_x reflects label uncertainties or, when labels are precise, a nominal value that sets the relative weighting in the loss function. Misspecified σ_x or σ_y leads to suboptimal correction: underestimating σ_x causes undercorrection, while overestimating it can cause overcorrection. A possible remedy is adding jitter terms to be optimized along with the other parameters.

LatentNN performs optimization rather than sampling. We obtain point estimates of $x_{\text{latent},i}$ and θ , but characterizing the full posterior over the high-dimensional space $(\theta, x_{\text{latent}})$ is challenging. Full Bayesian treatment of neural networks is already difficult, and adding the latent variables makes it more so. Approximate Bayesian approaches such as dropout (Hinton et al. 2012; Srivastava et al. 2014) or variational inference (Kingma & Welling 2013) may provide practical uncertainty estimates.

The increased number of parameters also leads to a more complex loss landscape. The optimization may converge to different solutions from different initializations, particularly at low SNR where the problem is ill-conditioned. Careful hyperparameter tuning, especially of the weight decay parameter as discussed in Section 5, is often necessary.

The computational cost scales as $\mathcal{O}(N \times p)$, introducing N additional parameters per input dimension. For full spectra with thousands of pixels and training sets of tens of thousands of stars, this becomes a consideration. Scalability requires care, particularly for high-dimensional inference problems. Possible strategies include using subsets of informative pixels. Dimensionality reduction before LatentNN is less straightforward because the noise properties in the reduced space may not be well characterized. Despite these limitations, the

from spectra to estimating photometric redshifts from imaging data. These applications often involve low-SNR observations where measurement uncertainties are non-negligible. When inputs are noisy, attenuation bias becomes a fundamental concern.

This linear regression model with coefficient slopes biased toward zero by a factor that scales with $\sigma_x/\sigma_{\text{noise}}$. Increasing model complexity does not fix this problem. **Crucial requirement: We need good estimates of the noise (σ_x, σ_y)! Garbage in, garbage out applies to noise estimates too.**

Applications involve high-SNR inputs where the effect is negligible. Astronomy is different, and as neural networks scale to larger datasets, as errors-in-variables become important. **It's computationally expensive! Scales with $N \times p$. Big data = big compute load.** Optimization is trickier with so many parameters to juggle.

By simultaneously learning from the data and noise, the model learns the correct input-output relationship rather than one diluted by measurement noise. **HEADS UP! It's not magic:**

We propose LatentNN as a practical implementation of this idea. **The problem: Neural nets inherit the 'squashing' bias from classic regression when inputs are noisy.**

Our solution: LatentNN! It learns the 'true' inputs as hidden variables, fixing the bias at the source.

The result: It works! We get unbiased, accurate models even with noisy data.

The big picture: This is a game-changer for astronomy, where noisy data is the norm.

WRAPPING IT UP:

Attenuation bias is a real danger in ML with noisy data, leading to wrong science. LatentNN is a rigorous, principled way to fix it by wedding modern deep learning with classic statistical ideas. It's not free, but for precision science, it's worth it!

Now go fix your models! - The Authors

Neural network modeling is a powerful tool, but it is not a panacea. Neural networks inherit the same limitations as classical statistical methods when those limitations arise from the data rather than the model. Attenuation bias is one such limitation. Recognizing when classical statistical insights apply to modern machine learning methods, and adapting solutions accordingly, is essential for reliable scientific inference in astronomical research.

CODE AVAILABILITY

We provide a reference implementation of LatentNN at <https://github.com/tingyuansen/LatentNN>. The code demonstrates the method for simple fully-connected architectures as presented in this paper. The approach is general and can be adapted to any differentiable architecture by treating the input as learnable parameters with an appropriate latent likelihood term.

Grab the code and start fixing your biases! →

ACKNOWLEDGMENTS

I thank David Weinberg for the encouragement that led to Paper I of this series. The original idea about

errors-in-variables stemmed from my teaching of Astron 5550 (Advanced Data Analysis), during which I decided to discuss attenuation bias as explored in Paper I. While preparing for that lecture, I was inspired to read more into the errors-in-variables literature, also suggested by an anonymous referee during Paper I. Further discussion with Jeffrey Newman and Brett Andrews during a colloquium at UPitt inspired the extension of this work to deep learning, which ultimately led to this paper. I also thank Hans-Walter Rix for useful discussion.

This research was supported by NSF Grant AST-2406729 and a Humboldt Research Award from the Alexander von Humboldt Foundation.

Claude 4.5 was used for assistance with code development and manuscript copy-editing. The scientific content, analysis, and conclusions are the author's, and all LLM-assisted changes were verified by the author.



Science is a team sport! Big thanks to everyone who bounced ideas around (probably over too much coffee).

A nod to our AI assistants tool 🤖

APPENDIX

APPENDIX: For the math folks, here's the derivation of the 'squashing'.

A. DERIVATION OF ATTENUATION BIAS

Consider a true linear relationship $y_{true} = \beta x_{true}$. Our observations are

$$\begin{aligned} x_{obs} &= x_{true} + \delta_x, \\ y_{obs} &= \beta x_{true} + \delta_y, \end{aligned}$$

where δ_x and δ_y are independent measurement errors with variances σ_x^2 and σ_y^2 .

The standard least squares estimator is $\hat{\beta} = \text{Cov}(y_{obs}, x_{obs}) / \text{Var}(x_{obs})$. Expanding

$$\begin{aligned} \text{Cov}(x_{obs}, y_{obs}) &= \text{Cov}(x_{true} + \delta_x, \beta x_{true} + \delta_y) \\ &= \beta \text{Var}(x_{true}) = \beta \sigma_{range}^2, \end{aligned} \tag{A3}$$

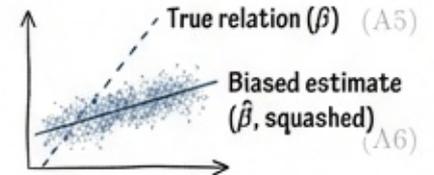
See? The math proves it. Noise in 'x' makes the estimated slope ($\hat{\beta}$) smaller than the true slope (β). That's attenuation bias in a nutshell. (A1)
(A2)
(A4)

since measurement errors are independent of true values and of each other. For the denominator:

$$\text{Var}(x_{obs}) = \text{Var}(x_{true}) + \text{Var}(\delta_x) = \sigma_{range}^2 + \sigma_x^2. \tag{A5}$$

The expected value of the standard estimator is therefore

$$\mathbb{E}[\hat{\beta}] = \beta \frac{\sigma_{range}^2}{\sigma_{range}^2 + \sigma_x^2} = \beta \cdot \lambda_\beta,$$



where $\lambda_\beta = 1 / (1 + (\sigma_x / \sigma_{range})^2) < 1$ is the attenuation factor. The estimated slope is systematically biased toward zero.

B. DERIVATION OF DEMING REGRESSION

And here's the classic statistical fix that inspired it all: Deming Regression.

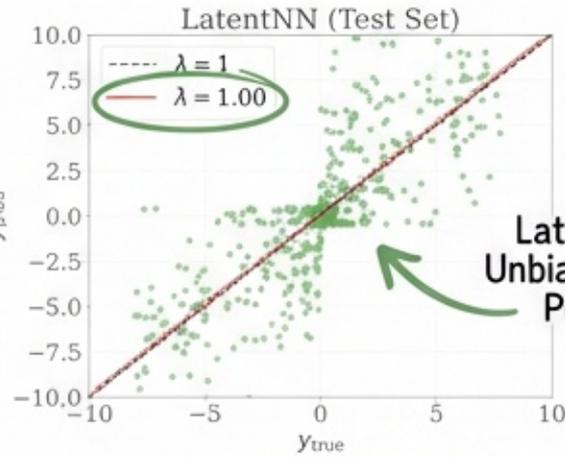
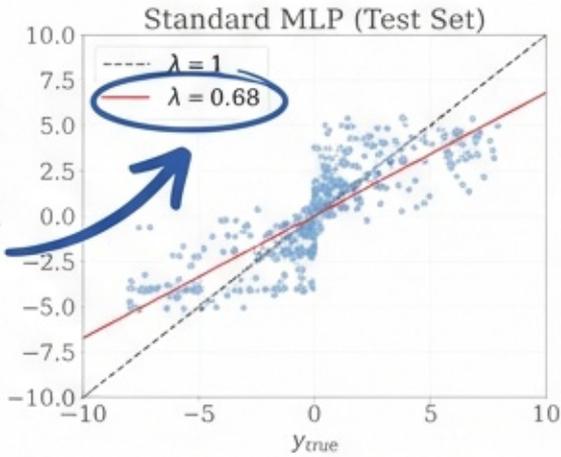
To correct attenuation bias, we treat the true values $\{x_{true,i}\}$ as latent variables to be estimated alongside β . The joint likelihood, conditioned on the true values, is

$$\mathcal{L} = \prod_{i=1}^N \mathcal{N}(x_{obs,i} | x_{true,i}, \sigma_x^2) \cdot \mathcal{N}(y_{obs,i} | \beta x_{true,i}, \sigma_y^2). \tag{B7}$$

The key idea: Treat the 'true' values ($x_{true,i}$) as unknowns (latent variables) and solve for them together with the slope.

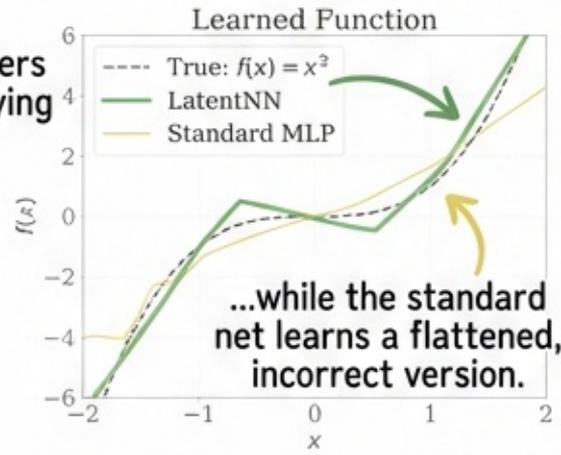
Seeing is believing: The non-linear challenge!

Standard Neural Net is STILL squashed by noise, even with a curvy function! Look at that bias! $\lambda < 1$.

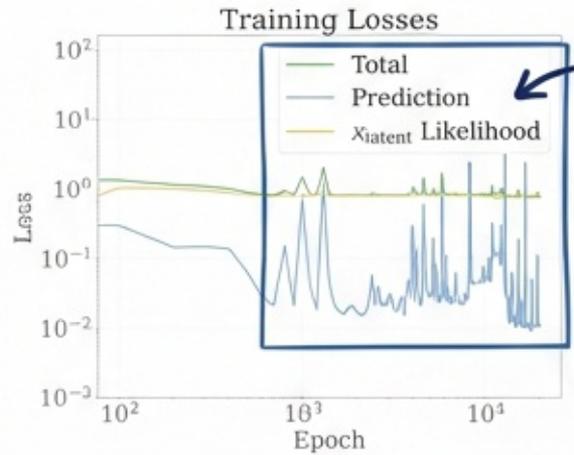


LatentNN nails it! Unbiased predictions. Perfect slope! $\lambda = 1$.

LatentNN discovers the TRUE underlying curve (x^3)...



...while the standard net learns a flattened, incorrect version.



The training "tug-of-war" in action! It's balancing fitting the data (Prediction loss) vs. finding plausible latent values (Likelihood loss).

Figure 7. LatentNN applied to a nonlinear function $f(x) = x^3$ at $SNR_x = 2$. Top: Predicted versus true y on the test set for standard MLP (left) and LatentNN (right). Bottom: Learned function (left) and training losses (right).

Back to the MATH ZONE:

How do we find those hidden "true" values?

Expanding the normal distributions and taking the derivative... This is the "cost function" we want to minimize. It penalizes estimates that are far from the observed noisy data (both x and y).

$$E(\{x_{true,i}\}, \beta) = \sum_{i=1}^N \left[\frac{(x_{obs,i} - x_{true,i})^2}{2\sigma_x^2} + \frac{(y_{obs,i} - \beta x_{true,i})^2}{2\sigma_y^2} \right] \quad (B6)$$

For a fixed β , each $x_{true,i}$ appears only in the i th term, so we can optimize them independently. Taking the derivative and setting it to zero:

$$\frac{\partial E}{\partial x_{true,i}} = -\frac{x_{obs,i} - x_{true,i}}{\sigma_x^2} - \frac{\beta(y_{obs,i} - \beta x_{true,i})}{\sigma_y^2} = 0$$

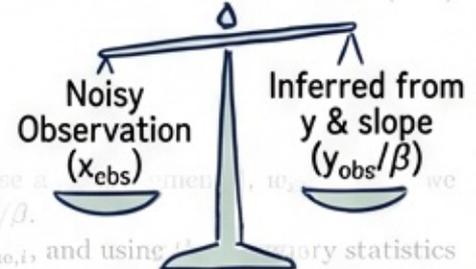
THE KEY INSIGHT!
The "true" x value is just a smart weighted average.

Rearranging and letting $r = \sigma_y^2/\sigma_x^2$ denote the error variance ratio:

$$x_{true,i} = \frac{r x_{obs,i} + \beta y_{obs,i}}{r + \beta^2}$$

This can be rewritten as a weighted average:

$$x_{true,i} = w_x \cdot x_{obs,i} + w_y \cdot \frac{y_{obs,i}}{\beta}$$



Weight depends on x-noise... vs y-noise. Clever!

where $w_x = r/(r + \beta^2)$ and $w_y = \beta^2/(r + \beta^2)$ sum to unity. When $r \ll 1$ (we trust x_{obs}), $w_x \approx 1$ and $w_y \approx 0$. Taking the derivative of E with respect to β , substituting the expression for $x_{true,i}$, and using the following statistics $S_{xx} = \sum_i (x_{obs,i} - \bar{x})^2$, $S_{yy} = \sum_i (y_{obs,i} - \bar{y})^2$, and $S_{xy} = \sum_i (x_{obs,i} - \bar{x})(y_{obs,i} - \bar{y})$, we obtain a quadratic equation:

$$S_{xy} \beta^2 - (S_{yy} - r S_{xx}) \beta - r S_{xy} = 0. \quad (B12)$$

Finally, we plug those "true" x 's back in to solve for the slope β (it gets a bit messy...).

Here's that messy result! The "true" slope β , finally revealed.

17

Applying the quadratic formula and taking the positive root:

Calculating...



$$\hat{\beta} = \frac{S_{yy} - rS_{xx} + \sqrt{(S_{yy} - rS_{xx})^2 + 4rS_{xy}^2}}{2S_{xy}} \quad (B13)$$

In the limit $\sigma_x \rightarrow 0$ (or equivalently $r \rightarrow \infty$), this reduces to $\hat{\beta} \rightarrow S_{xy}/S_{xx}$, the ordinary least squares estimator.

C. EXTENSION TO NONLINEAR FUNCTIONS

The main text focuses on linear function factor $\lambda_y = 1/(1 + (\sigma_s/\sigma_{\text{range}})^2)$. However, the errors-in-variables problem and the LatentNN correction are more general and apply to arbitrary nonlinear functions.

This text explains the cool plots on the LAST PAGE!

To demonstrate this, we test LatentNN on a cubic polynomial $f(x) = x^3$. Testing on a curvy function, not just a straight line. $\times 256$ hidden nnits) to ensure the network can represent the nonlinear function accurately. SNR_y = 2, a challenging regime where noise equals 50% of the signal range. The same weight decay hyperparameter search procedure from Section 5 is applied.

The standard MLP exhibits severe attenuation, with $\lambda_y \approx 0.65$ —even worse than the linear case at the same SNR, as derived in Paper 1. This occurs because measurement errors in x are amplified by the nonlinear function: an error δ_x produces a prediction error of approximately $f'(x)\delta_x$, which varies across the input domain. Near the edges of the training range where $|f'(x)| = 3x^2$ is largest, the effective noise is amplified.

LatentNN corrects this bias. The bottom-left panel shows the learned functions: the standard MLP learns a flattened version of the cubic, while LatentNN recovers the true function shape. This demonstrates that the latent variable framework extends naturally beyond linear regression to nonlinear function approximation.

LatentNN fixes it! It finds the true underlying curve. This method's just for simple the linear problems. It works for complex, real-world non-linear relationships too!

Andrae, R., Rix, H.-W., & Chandra, V. 2023a, ApJS, 267, 8, doi: 10.3847/1538-4365/acd53e

Andrae, R., Fouesneau, M., Sordo, R., et al. 2023b, A&A, 674, A27, doi: 10.1051/0004-6361/202243462

Bishop, C. M. 2006, Pattern Recognition and Machine Learning (Information Science and Statistics) (Berlin, Heidelberg: Springer-Verlag)

Blanco-Cuadros, S., Soubiran, C., Heiter, U., & Jofré, P. 2014, A&A, 569, A111, doi: 10.1051/0004-6361/201423845

Brown, T. B., Mann, B., Ryder, N., et al. 2020, arXiv e-prints, arXiv:2005.14165, doi: 10.48550/arXiv.2005.14165

Carroll, R. J., Ruppert, D., & Stefanski, L. A. 1985, Measurement Error in Non-linear Models (New York: Chapman and Hall)

Carroll, R. J., Ruppert, D., Stefanski, L. A., & Crainiceanu, C. M. 2006, Measurement error in nonlinear models: A modern perspective, Second edition (Chapman and Hall)

Chen, B., Orkney, M. D. A., Ting, Y.-S., & Hayden, M. R. 2025, arXiv e-prints, arXiv:2511.18901, doi: 10.48550/arXiv.2511.18901

Chen, B., Ting, Y.-S., & Hayden, M. 2024, PASA, 41, e063, doi: 10.1017/pasa.2024.56

Collister, A. A., & Lahav, O. 2004, PASP, 116, 345, doi: 10.1086/383254

Cranmer, K., Brehmer, J., & Louppe, G. 2020, Proceedings of the National Academy of Science, 117, 30055, doi: 10.1073/pnas.1912789117

Cybenko, G. 1989, Mathematics of Control, Signals, and Systems, 2, 303, doi: 10.1007/BF02551274

Deming, W. E. 1943, Statistical adjustment of data (Wiley)

Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. 2018, arXiv e-prints, arXiv:1810.04805, doi: 10.48550/arXiv.1810.04805

Fabbro, S., Venn, K. A., O'Brien, T., et al. 2018, MNRAS, 475, 2978, doi: 10.1093/mnras/stx3298

Fallows, C. P., & Sanders, J. L. 2024, MNRAS, 531, 2126, doi: 10.1093/mnras/stae1303

Frost, C., & Thompson, S. G. 2002, Journal of the Royal Statistical Society Series A: Statistics in Society, 163, 173, doi: 10.1111/1467-985X.00164

Fuller, W. A. 1987, Measurement Error Models (New York: Wiley)

García Pérez, A. E., Allende Prieto, C., Holtzman, J. A., et al. 2016, AJ, 151, 144, doi: 10.3847/0004-6250/151/6/144

Gelman, A., Carlin, J. B., Stern, H. S., et al. 2013, Bayesian data analysis (3rd ed.) (Chapman and Hall)

*The giants whose shoulders we stand on.
(Also, a great reading list for later.)*

And the list goes on... Seriously, science is a massive collaborative effort over decades!

Gheller, C., & Vazza, F. 2022, MNRAS, 509, 990, doi: 10.1093/mnras/stab3044

Golub, G. H., & Van Loan, C. F. 1980, SIAM Journal on Numerical Analysis, 17, 883

Guo, J., Bai, X., Deng, Y., et al. 2020, SoPh, 295, 3, doi: 10.1007/s11207-019-1573-9

Hattori, K. 2025, ApJ, 950, 90, doi: 10.3847/1538-4357/ad9686

He, K., Zhang, X., Ren, S., & Sun, J. 2016, in 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 1, doi: 10.1109/CVPR.2016.90

Hinton, G. E., Srivastava, N., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. R. 2012, arXiv e-prints, arXiv:1207.0580, doi: 10.48550/arXiv.1207.0580

Ho, A. Y. Q., Ness, M. K., Hogg, D. W., et al. 2017, ApJ, 836, 5, doi: 10.3847/1538-4337/836/1/5

Hoerl, A. E., & Kennard, R. W. 1970, Technometrics, 12, 55

Holtzman, J. A., Shetrone, M., Johnson, J. A., et al. 2015, AJ, 150, 148, doi: 10.1088/0004-6256/150/3/148

Hon, M., Huber, D., Kuzlewicz, J. S., et al. 2021, ApJ, 919, 131, doi: 10.3847/1538-4357/ac14b1

Hornik, K., Stinchcombe, M., & White, H. 1989, Neural Networks, 2, 359

Hoyle, B. 2016, Astronomy and Computing, 16, 34, doi: 10.1016/j.ascom.2016.03.006

Huertas-Company, M., & Lanusse, F. 2023, PASA, 20, e001, doi: 10.1017/pasa.2022.55

Kelly, B. C. 2007, ApJ, 665, 1489, doi: 10.1086/5177

Khalatyan, A., Anders, F., Chiappini, C., et al. 2024, MNRAS, 530, 691, A95, doi: 10.1051/0004-6361/202448100

Kingma, D. P., & Welling, M. 2013, arXiv e-prints, arXiv:1312.6114, doi: 10.48550/arXiv.1312.6114

Kurucz, R. L. 1970, SAO Special Report, 390

—. 1993, SYNTHSE spectrum synthesis program for stellar data

—. 2005, Memorie della Societa Astronomica Italiana Supplementi, 8, 14

—. 2017, ATLAS9: Model atmosphere program with opacity distribution functions, Astrophysics Source Code Library. <http://ascl.net/1710.017>

LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. 1998, Proceedings of the IEEE, 86, 2278

Leung, H. W., & Bovy, J. 2019, MNRAS, 483, 3253, doi: 10.1093/mnras/sty3217

Li, J., Wong, K. W. K., Hogg, D. W., Rix, H.-W., & Chandra, V. 2024, ApJS, 272, 2, doi: 10.3847/1538-4365/ad2b4d

Li, R., Napolitano, N. R., Feng, H., et al. 2022a, A&A, 666, A55, doi: 10.1051/0004-6361/202242100

Li, X., Wang, Z., Zeng, S., et al. 2022b, Research in Astronomy and Astrophysics, 22, 065018, doi: 10.1088/1674-4527/ac65e7

Lin, Q., Fouchez, D., Pasquet, J., et al. 2022, A&A, 662, A36, doi: 10.1051/0004-6361/202142751

Loredo, T. J. 2004, in American Institute of Physics Conference Series, Vol. 735, Bayesian Inference and Maximum Entropy Methods in Science and Engineering: 24th International Workshop on Bayesian Inference and Maximum Entropy Methods in Science and Engineering, ed. R. Fischer, R. Preuss, & U. V. Toussaint (AIP), 195-206, doi: 10.1063/1.1833214

Luo, A. L., Zhao, Y.-H., Zhao, G., et al. 2015, Research in Astronomy and Astrophysics, 15, 1095, doi: 10.1088/1674-4527/15/S/002

Majewski, S. R., Schiavon, R. P., Frinchaboy, P. M., et al. 2017, AJ, 154, 94, doi: 10.3847/1538-3881/aa784d

Ness, M., Hogg, D. W., Rix, H.-W., Ho, A. Y. Q., & Zuckerman, G. 2015, ApJ, 808, 16, doi: 10.1088/0004-6361/516/1/16

Odeh, M., Pennington, R. L., & Hogg, D. W. 1992, AJ, 103, 318, doi: 10.1086/112500

Opach, J., et al. 2023, arXiv e-prints, arXiv:2307.15850, doi: 10.48550/arXiv.2307.15850

Reid, B., & Dobbs, L. 2024, Astronomische Nachrichten, 345, e20240049, doi: 10.1002/asna.20240049

Reid, B., Pasquet, J., Chaumont, M., & Fouchez, D. 2019, A&A, 627, A21, doi: 10.1051/0004-6361/201834473

Rix, H.-W., Ting, Y.-S., Conroy, C., & Hogg, D. W. 2016, MNRAS, 461, L26, L25, doi: 10.3847/2041-8205/826/2/L25

Rix, H.-W., Chandra, V., Andrae, R., et al. 2022, ApJ, 941, 10, doi: 10.3847/1538-4357/ac9e01

Rubio-Delgado, D. E., Hinton, G. E., & Williams, R. J. 1986, Nature, 323, 833

Simonyan, K., & Zisserman, A. 2014, arXiv e-prints, arXiv:1409.1556, doi: 10.48550/arXiv.1409.1556

Spearman, C. 1904, The American Journal of Psychology, 15, 72, doi: 10.2307/1412159

Srivastava, N., Hinton, G. E., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. R. 2014, Journal of Machine Learning Research, 15, 1929

Storrie-Lombardi, M. C., Lahav, O., Sodre, Jr., L., & Storrie-Lombardi, L. J. 1992, MNRAS, 259, SP, doi: 10.1093/mnras/259.1.SP

Stranitz, I., Tkachenko, A., Gebruers, S., et al. 2022, AJ, 163, 236, doi: 10.3847/1538-3881/ac5f49

Szegedy, C., Liu, W., Jia, Y., et al. 2014, arXiv e-prints, arXiv:1409.1556, doi: 10.48550/arXiv.1409.1556

Talbot, C. H. 1938, Doklady, 4, 1035



See? Told you it's a great view from up here!

*Thanks for sticking with me to the end!
Now, go forth and stand on some shoulders of your own.
(And maybe read some of these papers!)*

A final few shout-outs to the community! Science is a team sport, after all.

- Ting, Y.-S., Conroy, C., Rix, H.-W., & Cargile, P. 2019, *ApJ*, 879, 69, doi: 10.3847/1538-4357/ab2331
- Ting, Y.-S., Rix, H.-W., Conroy, C., Ho, A. Y. Q., & Lin, J. 2017, *ApJL*, S49, L9, doi: 10.3847/2041-S213/aa921c
- Touvron, H., Lavril, T., Izacard, G., et al. 2023, arXiv e-prints, arXiv:2302.13971, doi: 10.48550/arXiv.2302.13971
- Van Huffel, S., & Vandewalle, J. 1991, *The total least squares problem: Computational aspects and analysis* (Society for Industrial and Applied Mathematics)
- Vojtekova, A., Lien, M., Valtchanov, I., et al. 2021, *MNRAS*, 503, 3204, doi: 10.1093/mnras/staa3567
- Vynatheya, P., Hamers, A. S., Mardling, R. A., et al. 2022, *MNRAS*, 510, 1000, doi: 10.1093/mnras/stab386
- Bellinger, E. P. 2022, *MNRAS*, 510, 1000, doi: 10.1093/mnras/stac238
- Hoyle, B. 2018, *Astronomy*, 10, 54, doi: 10.48550/arXiv:2203.00000
- Zuzuteosi, Coins, M. L., Roy, S., et al. 2014, *MNRAS*, 441, 8511, doi: 10.1093/mnras/stt2386
- Kelly, E. E., Conroy, T., Huang, Y., et al. 2014, *MNRAS*, 441, 8511, doi: 10.1093/mnras/stt2386
- Khalatyan, A., Xiang, N., L. & Aids, Bapriell 2014, *MNRAS*, 441, 691, doi: 10.1093/mnras/stt2386
- Kingma, D. P., & Welling, M. 2013, arXiv e-prints, arXiv:1312.6114, doi: 10.48550/arXiv.1312.6114
- Kurucz, R. L. 1970, *SAO Special Report*, 309
- . 1983, *SYNTHESIS spectrum synthesis program* data
- . 2005, *Memorie della Societa Astronomica Italiana Supplementi*, 37, 37, doi: 10.48550/arXiv.2005.14165
- Carroll, R. J., Ruppert, D., & Stefanski, L. A. 1988, *Measurement Error in Non-linear Models* (New York: Chapman and Hall)
- Carroll, R. J., Ruppert, D., Stefanski, L. A., & Crainiceann, C. M. 2006, *Measurement error in nonlinear models: A modern perspective*, Second edition (Chapman and Hall)
- Chen, B., Orkney, M. D. A., Ting, Y.-S., & Hayden, M. R. 2025, arXiv e-prints, arXiv:2511.15901, doi: 10.48550/arXiv.2511.15901
- Chen, B., Ting, Y.-S., & Hayden, M. 2024, *PASA*, 41, e063, doi: 10.1017/pasa.2024.55
- Wang, E. X., Nordlander, T., Asplund, M., et al. 2021, *MNRAS*, 500, 2159, doi: 10.1093/mnras/staa3351
- Winecki, D., & Kochanek, C. S. 2024, *ApJ*, 971, 61, doi: 10.3847/1538-4357/ad5a0b
- Wong, K. W. K., Ng, K. K. Y., & Berti, E. 2020, arXiv e-prints, arXiv:2007.10350, doi: 10.48550/arXiv.2007.10350
- Xiang, M., Ting, Y.-S., Rix, H.-W., et al. 2019, *ApJS*, 243, 34, doi: 10.3847/1538-4365/ab3364
- Xiang, M., Rix, H.-W., Ting, Y.-S., et al. 2022, *A&A*, 662, A66, doi: 10.1051/0004-6361/202141570
- Yang, L., Yuan, H., Huang, B., et al. 2025, *ApJS*, 279, 7, doi: 10.3847/1538-4365/add5e3
- Zhang, M., Xiang, M., Ting, Y.-S., et al. 2024, *ApJS*, 273, 19, doi: 10.3847/1538-4363/ad51dd
- Zhang, X., Green, G. M., & Rix, H.-W. 2023, *MNRAS*, S24, 1855, doi: 10.1093/mnras/stad1941
- Zhang, X., Green, G. M., & Rix, H.-W. 2022, *MNRAS*, 510, 1000, doi: 10.1093/mnras/stac786

THE END.



Thanks for reading!
Now, where's that coffee?

PS: Don't forget to check out the code repository mentioned earlier!