

## A APPENDIX

### B PROOF OF EQUATION (2)

As a reminder, we consider  $C$  class labels and denote by  $\Delta_C$  the  $C$ -dimensional simplex. We define the set of distributions  $\mathcal{Q}$  over the  $C$  class labels by

$$\mathcal{Q} = \{q(y|\cdot) \mid \forall \mathbf{x} \in \mathcal{X}, q(y|\mathbf{x}) \in \Delta_C\}.$$

Consider the optimization problem

$$\min_{q \in \mathcal{Q}} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} [D_{KL}(p(y|\mathbf{x}) \parallel q(y|\mathbf{x}))] \quad (5)$$

whose solution is straightforwardly given by the marginal distribution  $\mathbf{x} \mapsto q^*(y|\mathbf{x}) = p(y|\mathbf{x})$ . We recall that the KL  $D_{KL}(p(y|\mathbf{x}) \parallel q(y|\mathbf{x}))$  is defined by

$$D_{KL}(p(y|\mathbf{x}) \parallel q(y|\mathbf{x})) = \sum_{j=1}^C p_j(y|\mathbf{x}) \log \left( \frac{p_j(y|\mathbf{x})}{q_j(y|\mathbf{x})} \right). \quad (6)$$

For any  $\mathbf{x}$  and  $j \leq C$ , we can rewrite the terms of the sum

$$p_j(y|\mathbf{x}) \log \left( \frac{p_j(y|\mathbf{x})}{q_j(y|\mathbf{x})} \right)$$

as

$$\mathbb{E}_{\mathbf{a}|\mathbf{x} \sim p(\mathbf{a}|\mathbf{x})} \left[ p_j(y|\mathbf{x}, \mathbf{a}) \log \left( \frac{p_j(y|\mathbf{x})}{q_j(y|\mathbf{x})} \right) \right]$$

where we have used (i) the fact that  $\log \left( \frac{p_j(y|\mathbf{x})}{q_j(y|\mathbf{x})} \right)$  does not depend on  $\mathbf{a}$  and (ii) the definition of the marginal distribution

$$\begin{aligned} p_j(y|\mathbf{x}) &= \int p_j(y|\mathbf{x}, \mathbf{a}) p(\mathbf{a}|\mathbf{x}) d\mathbf{a} \\ &= \mathbb{E}_{\mathbf{a}|\mathbf{x} \sim p(\mathbf{a}|\mathbf{x})} [p_j(y|\mathbf{x}, \mathbf{a})]. \end{aligned}$$

Multiplying and dividing in the argument of the log by  $p_j(y|\mathbf{x}, \mathbf{a})$ , we obtain

$$\mathbb{E}_{\mathbf{a}|\mathbf{x} \sim p(\mathbf{a}|\mathbf{x})} \left[ p_j(y|\mathbf{x}, \mathbf{a}) \log \left( \frac{p_j(y|\mathbf{x}, \mathbf{a})}{q_j(y|\mathbf{x})} \frac{p_j(y|\mathbf{x})}{p_j(y|\mathbf{x}, \mathbf{a})} \right) \right].$$

Summing over  $j \in \{1, \dots, C\}$  to reconstruct the KL term (6), this leads to, for any  $\mathbf{x}$ ,

$$\begin{aligned} D_{KL}(p(y|\mathbf{x}) \parallel q(y|\mathbf{x})) &= \mathbb{E}_{\mathbf{a}|\mathbf{x}} [D_{KL}(p(y|\mathbf{x}, \mathbf{a}) \parallel q(y|\mathbf{x}))] \\ &\quad - \mathbb{E}_{\mathbf{a}|\mathbf{x}} [D_{KL}(p(y|\mathbf{x}, \mathbf{a}) \parallel p(y|\mathbf{x}))]. \end{aligned}$$

Since the second term above *does not depend on*  $q$ , minimizing (5) is equivalent to minimizing

$$\begin{aligned} &\min_{q \in \mathcal{Q}} \mathbb{E}_{\mathbf{x}} [\mathbb{E}_{\mathbf{a}|\mathbf{x}} [D_{KL}(p(y|\mathbf{x}, \mathbf{a}) \parallel q(y|\mathbf{x}))]] \\ &= \min_{q \in \mathcal{Q}} \mathbb{E}_{(\mathbf{x}, \mathbf{a}) \sim p(\mathbf{x}, \mathbf{a})} [D_{KL}(p(y|\mathbf{x}, \mathbf{a}) \parallel q(y|\mathbf{x}))] \end{aligned}$$

which is equal to (2) and which is, analogously to (5), minimized by the marginal distribution  $\mathbf{x} \mapsto q^*(y|\mathbf{x}) = p(y|\mathbf{x})$ .

### C HETEROSCEDASTIC MOTIVATION

We consider a simplified special case of our framework in which the conditional model  $p(y|\mathbf{x}, \mathbf{a})$  is *homoscedastic* but the optimal variational distribution in the sense of Eq. 2 is *heteroscedastic*. This motivates **Het-TRAM**, in which the variational approximations  $q(y|\mathbf{x})$  and  $q(y|\mathbf{x}, \mathbf{a})$  are heteroscedastic.

Suppose we have a regression dataset constructed from labels assigned by  $M$  annotators. Each annotator has their own homoscedastic Gaussian model  $p(y|\mathbf{x}, a = m) = \mathcal{N}(\mu_{\theta_m}(\mathbf{x}), 1)$ . Here the

Table 4: Pre-trained models used to re-label ImageNet ILSVRC12 training set and their accuracy on that training set.

Model	Training set accuracy
ResNet50V2	0.70086
ResNet101V2	0.72346
ResNet152V2	0.72738
DenseNet121	0.74782
DenseNet169	0.76184
DenseNet201	0.77344
InceptionResNetV2	0.8049
InceptionV3	0.77994
MobileNet	0.70594
MobileNetV2	0.71458
MobileNetV3Large	0.75622
MobileNetV3Small	0.68158
NASNetMobile	0.74302
VGG16	0.71178
VGG19	0.71156
Xception	0.79076

PI is a single discrete Categorical feature representing the annotator ID which takes one of  $M$  values with equal probability,  $a \sim \text{Cat}(\frac{1}{M})$ .

The marginal  $p(y|\mathbf{x})$  is a Gaussian Mixture Model. We choose our variational family to be the Gaussian distribution,  $q(y|\mathbf{x}) = \mathcal{N}(\mu(\mathbf{x}), \sigma^2(\mathbf{x}))$ . The values of  $\mu$  and  $\sigma^2$  that minimize Eq. 2 are:  $\mu_*(\mathbf{x}) = \frac{1}{M} \sum_m \mu_{\theta_m}(\mathbf{x})$  and  $\sigma_*^2(\mathbf{x}) = (M - \mu_*(\mathbf{x})) + \frac{1}{M} \sum_m \mu_{\theta_m}^2(\mathbf{x})$  (Lakshminarayanan et al., 2017). Crucially note that despite the conditional distribution being homoscedastic, the best variational distribution is heteroscedastic as the variance depends on the location in  $\mathcal{X}$  space.

## D EXPERIMENTAL DETAILS

### D.1 DATA GENERATION PROCESS

**CIFAR-10H.** We use the CIFAR-10 image as the non-privileged information  $\mathbf{x}$ . The annotator ID, the number of prior annotations the annotator has provided and the reaction time in milliseconds of the annotator, are used as privileged information  $\mathbf{a}$ . For feature pre-processing the annotator ID is one-hot encoded. The number of prior annotations and the reaction time are independently divided into 10 equally sized quantiles and the quantile ID is one-hot encoded. The image is pre-processed according to the standard MobileNet pre-processing (Howard et al., 2017).

As CIFAR-10H has on average more than 50 annotations per image and the labels are not particularly noisy. We subsample the CIFAR-10H labels by the following procedure. We keep all labels by the 41 annotators that agree with the true CIFAR-10 label less than 85% of the time. We then select a further 41 annotators from the remaining annotators. The average agreement of the bad annotators with the CIFAR-10 label is 63.3%, in the full subset of labels: 79.2% and in the full CIFAR-10H dataset: 94.9%. The subsampling procedure leaves 16,400 labels from 82 annotators while the full CIFAR-10H dataset has 514,200 labels from 2,571 annotators.

**ImageNet.** The annotator features are the model ID used to re-label  $\mathbf{x}$ , which is one-hot encoded and the probability of that label being sampled. See main paper for details on the sampling procedure and see Table 4 for the list of models used and their accuracy on the ImageNet training set. The pre-trained models are downloaded from tf.keras.applications<sup>2</sup>.

<sup>2</sup>[https://www.tensorflow.org/api\\_docs/python/tf/keras/applications](https://www.tensorflow.org/api_docs/python/tf/keras/applications)

## D.2 HYPERPARAMETERS

**CIFAR-10H.** For all methods  $\phi(\mathbf{x})$  (or equivalent) is a MobileNet (Howard et al., 2017) pre-trained on ImageNet ILSVRC12, followed by a global average pooling layer and a Dense + ReLU layer with 64 units.  $\psi(\mathbf{x}, \mathbf{a})$  is a two-layer MLP with 64 units per layer and ReLU activation. The first layer takes only  $\mathbf{a}$  as an input, while the second layer takes the output of the first layer concatenated with  $\phi(\mathbf{x})$  as input.

All models are trained for 20 epochs with the Adam optimizer with base learning rate = 0.001,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\epsilon = 1e - 07$ . All models are trained with L2 weight regularization with weighting  $1e - 3$ .

Heteroscedastic models are trained using the method of Collier et al. (2021) with 4 factors for the low-rank covariance matrix approximation and a softmax temperature parameter of  $\tau = 3.0$ . Distilled models are also trained with a softmax temperature of  $\tau = 3.0$  to smooth the teacher labels and with the distillation hyperparameter  $\lambda = 0.5$  which weights the losses from the soft teacher labels and the true labels. A grid search over  $\tau \in \{1.0, 2.0, 3.0, 4.0\}$  and  $\lambda \in \{0.0, 0.25, 0.5, 0.75, 1.0\}$  was conducted.

**ImageNet.** For all methods  $\phi(\mathbf{x})$  (or equivalent) is a randomly initialized ResNet-50 (He et al., 2016) with the output layer removed.  $\psi(\mathbf{x}, \mathbf{a})$  is a two-layer MLP with 128 units per layer and ReLU activation, the output of this MLP is concatenated with  $\phi(\mathbf{x})$  and then passed to the output layer. The first layer of the  $\psi(\mathbf{x}, \mathbf{a})$  MLP takes only  $\mathbf{a}$  as an input, while the second layer takes the output of the first layer concatenated with  $\phi(\mathbf{x})$  as input.

All but Het-TRAM models are trained for 90 epochs with the SGD optimizer with base learning rate = 0.1, decayed by a factor of 10 after 30, 60 and 80 epochs. Following Collier et al. (2021), Het-TRAM is trained for 270 epochs with the same initial learning rate and learning rate decay at 90, 180 and 240 epochs. All models are trained with L2 weight regularization with weighting  $1e - 4$ .

Heteroscedastic models use 15 factors for the low-rank covariance matrix approximation and a softmax temperature parameter of  $\tau = 1.5$ . Distilled models are trained with a softmax temperature of  $\tau = 3.0$  and with the distillation hyperparameter  $\lambda = 0.5$ . A grid search over  $\tau \in \{1.0, 2.0, 3.0, 4.0\}$  and  $\lambda \in \{0.0, 0.25, 0.5, 0.75, 1.0\}$  was conducted.

## E RISK ANALYSIS

**Generative model and notations.** We assume the following

- $\mathbf{a} \in \mathbb{R}^m, \mathbf{x} \in \mathbb{R}^d$ ,
- $\mathbf{a} \sim p(\mathbf{a}|\mathbf{x}) = \mathcal{N}(\mu(\mathbf{x})|\Sigma(\mathbf{x}))$  for some mean and covariance dependent on  $\mathbf{x}$ ,
- $\mathbf{y} = \mathbf{x}^\top \mathbf{w}^* + \mathbf{a}^\top \mathbf{v}^* + \varepsilon$  with  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ .

When considering  $n$  observations from this generative model, we use the matrix representations  $\mathbf{y} \in \mathbb{R}^n, \mathbf{X} \in \mathbb{R}^{n \times d}, \mathbf{A} \in \mathbb{R}^{n \times m}$  and  $\varepsilon \in \mathbb{R}^n$ . We also write the zero-mean Gaussian vector

$$\mathbf{z} = (\mathbf{A} - \mu(\mathbf{X}))\mathbf{v}^* + \varepsilon \in \mathbb{R}^n \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I} + \mathbf{\Lambda})$$

where we have defined the diagonal covariance

$$\mathbf{\Lambda} = \mathbf{\Lambda}(\mathbf{v}^*, \mathbf{X}) = \text{Diag}(\{(\mathbf{v}^*)^\top \Sigma(\mathbf{x}_i) \mathbf{v}^*\}_{i=1}^n) \in \mathbb{R}^{n \times n}.$$

We list below some notation that we will repeatedly use

- The orthogonal projector associated with  $\mathbf{X}$ :

$$\mathbf{\Pi}_x = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \in \mathbb{R}^{n \times n}.$$

- Similarly, the orthogonal projector associated with  $\mathbf{A}$ :

$$\mathbf{\Pi}_a = \mathbf{A}(\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \in \mathbb{R}^{n \times n}.$$

- The projections  $\mathbf{X}_{a\perp} = (\mathbf{I} - \mathbf{\Pi}_a)\mathbf{X}$  and  $\mathbf{A}_{x\perp} = (\mathbf{I} - \mathbf{\Pi}_x)\mathbf{A}$ .
- The matrices:  $\mathbf{H} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \in \mathbb{R}^{d \times n}$  and  $\mathbf{G} = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \in \mathbb{R}^{m \times n}$ .
- The matrices above when restricted to the projections of  $\mathbf{X}$  and  $\mathbf{A}$  respectively, that is,  

$$\mathbf{H}_{a\perp} = (\mathbf{X}_{a\perp}^\top \mathbf{X}_{a\perp})^{-1} \mathbf{X}_{a\perp}^\top \in \mathbb{R}^{d \times n} \text{ and } \mathbf{G}_{x\perp} = (\mathbf{A}_{x\perp}^\top \mathbf{A}_{x\perp})^{-1} \mathbf{A}_{x\perp}^\top \in \mathbb{R}^{m \times n}.$$

### E.1 DEFINITION OF THE RISK

We will compare different estimators based on their different *risks*. We focus on the *fixed* design analysis (Bach, 2021), i.e., we study the errors only due to resampling the noise  $\varepsilon$  and the feature  $\mathbf{a}$ .

Given a predictor  $\tau(\mathbf{X})$  based on the training quantities  $(\mathbf{X}, \mathbf{A}, \varepsilon)$ , we consider  $\mathbf{y}' = \mathbf{X}\mathbf{w}^* + \mathbf{A}'\mathbf{v}^* + \varepsilon'$  (where the prime is to stress the difference with the training quantities without prime) and define the risk of  $\tau$  as

$$\mathcal{R}(\tau(\mathbf{X})) = \mathbb{E}_{\varepsilon' \sim p(\varepsilon'), \mathbf{a}' \sim p(\mathbf{a}'|\mathbf{x})} \left\{ \frac{1}{n} \|\mathbf{y}' - \tau(\mathbf{X})\|^2 \right\}. \quad (7)$$

Expanding the square with  $\mathbf{y}' - \tau(\mathbf{X}) = \mathbf{X}\mathbf{w}^* - \tau(\mathbf{X}) + \mu(\mathbf{X})\mathbf{v}^* + \mathbf{z}'$ , we obtain the expression

$$\mathcal{R}(\tau(\mathbf{X})) = \frac{1}{n} \|\mathbf{X}\mathbf{w}^* - \tau(\mathbf{X}) + \mu(\mathbf{X})\mathbf{v}^*\|^2 + \frac{1}{n} \text{tr}(\sigma^2 \mathbf{I} + \mathbf{\Lambda}). \quad (8)$$

Following common practices (Bach, 2021), to assess the risk, we finally take a second expectation  $\mathbb{E}_{\varepsilon \sim p(\varepsilon), \mathbf{a} \sim p(\mathbf{a}|\mathbf{x})} [\mathcal{R}(\tau(\mathbf{X}))]$  with respect to the training quantities  $(\mathbf{A}, \varepsilon)$ .

Since we will mostly consider differences of risks, we omit the variance term  $\frac{1}{n} \text{tr}(\sigma^2 \mathbf{I} + \mathbf{\Lambda})$  in the equations below.

### E.2 CAPTURING THE BENEFIT OF PI WITHOUT MARGINALIZATION

We first describe when, in absence of any marginalization, ordinary least squares ignoring PI is worse than ordinary least squares using PI with mean imputation at prediction time.

**Proposition E.1.** *Assume that  $\mathbf{X}^\top \mathbf{X}$  is invertible. Moreover, assume that  $\mathbf{A}^\top \mathbf{A}$  and  $[\mathbf{X}, \mathbf{A}]^\top [\mathbf{X}, \mathbf{A}]$  are almost surely invertible. We have that*

$$\mathbb{E}[\mathcal{R}(\tau_{\text{NO-PI}}(\mathbf{X}))] > \mathbb{E}[\mathcal{R}(\tau_{\text{PI}}(\mathbf{X}))]$$

if and only if

$$\frac{1}{n} \|( \mathbf{I} - \mathbf{\Pi}_x ) \mu(\mathbf{X}) \mathbf{v}^*\|^2 + \frac{\sigma^2 d}{n} + \frac{1}{n} \text{tr}(\mathbf{\Pi}_x \mathbf{\Lambda}) > \frac{\sigma^2}{n} \mathbb{E}[\|\mathbf{K}\|^2]$$

with  $\mathbf{K} = \mathbf{X}\mathbf{H}_{a\perp} + \mu(\mathbf{X})\mathbf{G}_{x\perp}$ . When  $m = 1$  (i.e.,  $\mathbf{A}$  is a column vector), a sufficient condition is

$$\frac{1}{n} \|( \mathbf{I} - \mathbf{\Pi}_x ) \mu(\mathbf{X}) \mathbf{v}^*\|^2 + \frac{1}{n} \text{tr}(\mathbf{\Pi}_x \mathbf{\Lambda}) > 2 \mathbb{E} \left[ \frac{\|\mathbf{\Pi}_x \mathbf{A}\|^2 + \|\mu(\mathbf{X})\|^2}{\|( \mathbf{I} - \mathbf{\Pi}_x ) \mathbf{A}\|^2} \right] + \frac{\sigma^2 d}{n}.$$

We provide the details of the derivation of the risk for  $\tau_{\text{NO-PI}}$  and  $\tau_{\text{PI}}$  in Section E.2.1 and Section E.2.2 respectively. Moreover, the second part of the proposition follows from an application of Lemma E.5.

#### E.2.1 ORDINARY LEAST SQUARES (NO MARGINALIZATION)

The solution of

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$$

is given by  $\hat{\mathbf{w}}_0 = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \mathbf{H}\mathbf{y}$ . The corresponding predictions are

$$\tau_{\text{NO-PI}}(\mathbf{X}) = \mathbf{X}\hat{\mathbf{w}}_0 = \mathbf{\Pi}_x \mathbf{y} = \mathbf{X}\mathbf{w}^* + \mathbf{\Pi}_x \mu(\mathbf{X}) \mathbf{v}^* + \mathbf{\Pi}_x \mathbf{z}.$$

Plugging into (8), we obtain

$$\mathcal{R}(\tau_{\text{NO-PI}}(\mathbf{X})) = \frac{1}{n} \|( \mathbf{I} - \mathbf{\Pi}_x ) \mu(\mathbf{X}) \mathbf{v}^* - \mathbf{\Pi}_x \mathbf{z}\|^2.$$

Expanding the square and using that  $\text{tr}(\mathbf{\Pi}_x) = d$ , the final risk expression is

$$\begin{aligned} \mathbb{E}[\mathcal{R}(\tau_{\text{NO-PI}}(\mathbf{X}))] &= \frac{1}{n} \|( \mathbf{I} - \mathbf{\Pi}_x ) \mu(\mathbf{X}) \mathbf{v}^*\|^2 + \frac{1}{n} \mathbb{E}[\|\mathbf{\Pi}_x \mathbf{z}\|^2] \\ &= \frac{1}{n} \|( \mathbf{I} - \mathbf{\Pi}_x ) \mu(\mathbf{X}) \mathbf{v}^*\|^2 + \frac{\sigma^2 d}{n} + \frac{1}{n} \text{tr}(\mathbf{\Pi}_x \mathbf{\Lambda}). \end{aligned} \quad (9)$$

### E.2.2 ORDINARY LEAST SQUARES WITH PI AND MEAN IMPUTATION (NO MARGINALIZATION)

We focus on the solution of

$$\min_{\mathbf{w}, \mathbf{v}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w} - \mathbf{A}\mathbf{v}\|^2$$

to construct an estimator. Using Lemma E.3 we have

$$\hat{\mathbf{w}}_1 = \mathbf{H}_{a\perp} \mathbf{y} \quad \text{and} \quad \hat{\mathbf{v}}_1 = \mathbf{G}_{x\perp} \mathbf{y}.$$

Using Lemma E.4 we can simplify

$$\hat{\mathbf{w}}_1 = \mathbf{H}_{a\perp} \mathbf{y} = \mathbf{w}^* + \mathbf{0} + \mathbf{H}_{a\perp} \varepsilon$$

and

$$\hat{\mathbf{v}}_1 = \mathbf{G}_{x\perp} \mathbf{y} = \mathbf{0} + \mathbf{v}^* + \mathbf{G}_{x\perp} \varepsilon.$$

Since  $\mathbf{A}$  is not available at prediction time, we impute it instead with its mean  $\mu(\mathbf{X})$ , which is assumed to be perfectly known. This leads to

$$\tau_{\text{PI}}(\mathbf{X}) = \mathbf{X}\hat{\mathbf{w}}_1 + \mu(\mathbf{X})\hat{\mathbf{v}}_1 = \mathbf{X}\mathbf{w}^* + \mu(\mathbf{X})\mathbf{v}^* + \mathbf{K}\varepsilon,$$

with

$$\mathbf{K} = \mathbf{X}\mathbf{H}_{a\perp} + \mu(\mathbf{X})\mathbf{G}_{x\perp}.$$

Plugging into (8) and taking the expectation, we obtain

$$\begin{aligned} \mathbb{E}[\mathcal{R}(\tau_{\text{PI}}(\mathbf{X}))] &= \frac{1}{n} \|\mathbf{0}\|^2 + \frac{1}{n} \mathbb{E}[\|\mathbf{K}\varepsilon\|^2] \\ &= \frac{\sigma^2}{n} \mathbb{E}[\|\mathbf{K}\|^2]. \end{aligned} \quad (10)$$

### E.3 CAPTURING THE BENEFIT OF PI WITH MARGINALIZATION

We then describe when, with marginalization, ordinary least squares ignoring PI is worse than ordinary least squares using PI.

**Proposition E.2.** *Assume that  $\mathbf{X}^\top \mathbf{X}$  is invertible. Moreover, assume that  $\mathbf{A}^\top \mathbf{A}$  and  $[\mathbf{X}, \mathbf{A}]^\top [\mathbf{X}, \mathbf{A}]$  are almost surely invertible. We have that*

$$\mathbb{E}[\mathcal{R}(\tau_{\text{marg. NO-PI}}(\mathbf{X}))] > \mathbb{E}[\mathcal{R}(\tau_{\text{marg. PI}}(\mathbf{X}))]$$

if and only if

$$\frac{1}{n} \|(\mathbf{I} - \mathbf{\Pi}_x) \mu(\mathbf{X}) \mathbf{v}^*\|^2 + \frac{\sigma^2 d}{n} > \frac{\sigma^2}{n} \|\mathbb{E}[\mathbf{L}]\|^2$$

with  $\mathbf{L} = \mathbf{X}\mathbf{H}_{a\perp} + \mathbf{A}\mathbf{G}_{x\perp}$ . When  $m = 1$  (i.e.,  $\mathbf{A}$  is a column vector), a sufficient condition is

$$\frac{1}{n} \|(\mathbf{I} - \mathbf{\Pi}_x) \mu(\mathbf{X}) \mathbf{v}^*\|^2 > 2 \mathbb{E} \left[ \frac{\|\mathbf{\Pi}_x \mathbf{A}\|^2 + \|\mathbf{A}\|^2}{\|(\mathbf{I} - \mathbf{\Pi}_x) \mathbf{A}\|^2} \right] + \frac{\sigma^2 d}{n}.$$

We provide the details of the derivation of the risk for  $\tau_{\text{marg. NO-PI}}$  and  $\tau_{\text{marg. PI}}$  in Section E.3.1 and Section E.3.2 respectively. Moreover, the second part of the proposition follows from an application of Lemma E.5.

#### E.3.1 ORDINARY LEAST SQUARES (WITH MARGINALIZATION)

Restarting from Section E.2.1, we consider the predictions marginalized with respect to  $\mathbf{A}$ . We have

$$\tau_{\text{marg. NO-PI}}(\mathbf{X}) = \mathbb{E}_{\mathbf{a} \sim p(\mathbf{a}|\mathbf{x})} [\mathbf{X}\hat{\mathbf{w}}_0] = \mathbf{X}\mathbf{w}^* + \mathbf{\Pi}_x \mu(\mathbf{X}) \mathbf{v}^* + \mathbf{\Pi}_x \varepsilon.$$

Plugging into (8), we obtain

$$\mathcal{R}(\tau_{\text{marg. NO-PI}}(\mathbf{X})) = \frac{1}{n} \|(\mathbf{I} - \mathbf{\Pi}_x) \mu(\mathbf{X}) \mathbf{v}^* - \mathbf{\Pi}_x \varepsilon\|^2.$$

Expanding the square and using that  $\text{tr}(\mathbf{\Pi}_x) = d$ , the final risk expression is

$$\begin{aligned} \mathbb{E}[\mathcal{R}(\tau_{\text{marg. NO-PI}}(\mathbf{X}))] &= \frac{1}{n} \|(\mathbf{I} - \mathbf{\Pi}_x) \mu(\mathbf{X}) \mathbf{v}^*\|^2 + \frac{1}{n} \mathbb{E}[\|\mathbf{\Pi}_x \varepsilon\|^2] \\ &= \frac{1}{n} \|(\mathbf{I} - \mathbf{\Pi}_x) \mu(\mathbf{X}) \mathbf{v}^*\|^2 + \frac{\sigma^2 d}{n}. \end{aligned} \quad (11)$$

### E.3.2 ORDINARY LEAST SQUARES WITH PI AND MARGINALIZATION

Restarting from Section E.2.2 we consider the predictions marginalized with respect to  $\mathbf{A}$ . In particular, we do not impute  $\mathbf{A}$  by its mean but rather directly take the expectation over  $\mathbf{A}$ . We have

$$\tau_{\text{marg. PI}}(\mathbf{X}) = \mathbb{E}_{\mathbf{a} \sim p(\mathbf{a}|\mathbf{x})}[\mathbf{X}\hat{\mathbf{w}}_1 + \mathbf{A}\hat{\mathbf{v}}_1] = \mathbf{X}\mathbf{w}^* + \mu(\mathbf{X})\mathbf{v}^* + \mathbb{E}_{\mathbf{a} \sim p(\mathbf{a}|\mathbf{x})}[\mathbf{L}]\varepsilon,$$

with

$$\mathbf{L} = \mathbf{X}\mathbf{H}_{a\perp} + \mathbf{A}\mathbf{G}_{x\perp}.$$

Plugging into (8) and taking the expectation, we obtain

$$\begin{aligned} \mathbb{E}[\mathcal{R}(\tau_{\text{marg. PI}}(\mathbf{X}))] &= \frac{1}{n}\|\mathbf{0}\|^2 + \frac{1}{n}\mathbb{E}[\|\mathbb{E}_{\mathbf{a} \sim p(\mathbf{a}|\mathbf{x})}[\mathbf{L}]\varepsilon\|^2] \\ &= \frac{\sigma^2}{n}\|\mathbb{E}_{\mathbf{a} \sim p(\mathbf{a}|\mathbf{x})}[\mathbf{L}]\|^2. \end{aligned} \quad (12)$$

### E.4 TECHNICAL LEMMAS

**Lemma E.3.** Assume that both  $\mathbf{X}^\top \mathbf{X}$  and  $\mathbf{A}^\top \mathbf{A}$  are invertible. Moreover, assume that both  $\mathbf{X}_{a\perp}^\top \mathbf{X}_{a\perp}$  and  $\mathbf{A}_{x\perp}^\top \mathbf{A}_{x\perp}$  are invertible.

We can write the solution of

$$\min_{\mathbf{w}, \mathbf{v}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w} - \mathbf{A}\mathbf{v}\|^2$$

as

$$\hat{\mathbf{w}} = \mathbf{H}_{a\perp} \mathbf{y} \quad \text{and} \quad \hat{\mathbf{v}} = \mathbf{G}_{x\perp} \mathbf{y}.$$

*Proof.* The proof follows by applying inversion formula for the block matrix

$$\mathbf{Q} = \begin{bmatrix} \mathbf{X}^\top \mathbf{X} & \mathbf{X}^\top \mathbf{A} \\ \mathbf{A}^\top \mathbf{X} & \mathbf{A}^\top \mathbf{A} \end{bmatrix}$$

where  $\mathbf{X}_{a\perp}^\top \mathbf{X}_{a\perp}$  and  $\mathbf{A}_{x\perp}^\top \mathbf{A}_{x\perp}$  are the two Schur complements of  $\mathbf{X}^\top \mathbf{X}$  and  $\mathbf{A}^\top \mathbf{A}$ . Under the assumptions of the lemma, the matrix is  $\mathbf{Q}$  is invertible.  $\square$

**Lemma E.4.** We have the following properties

- $\mathbf{H}_{a\perp} \mathbf{X} = (\mathbf{X}_{a\perp}^\top \mathbf{X}_{a\perp})^{-1} \mathbf{X}^\top (\mathbf{I} - \mathbf{\Pi}_a) \mathbf{X} = (\mathbf{X}_{a\perp}^\top \mathbf{X}_{a\perp})^{-1} (\mathbf{X}_{a\perp}^\top \mathbf{X}_{a\perp}) = \mathbf{I},$
- $\mathbf{H}_{a\perp} \mathbf{A} = (\mathbf{X}_{a\perp}^\top \mathbf{X}_{a\perp})^{-1} \mathbf{X}^\top (\mathbf{I} - \mathbf{\Pi}_a) \mathbf{A} = \mathbf{0}.$

Conversely, we have

- $\mathbf{G}_{x\perp} \mathbf{A} = (\mathbf{A}_{x\perp}^\top \mathbf{A}_{x\perp})^{-1} \mathbf{A}^\top (\mathbf{I} - \mathbf{\Pi}_x) \mathbf{A} = (\mathbf{A}_{x\perp}^\top \mathbf{A}_{x\perp})^{-1} (\mathbf{A}_{x\perp}^\top \mathbf{A}_{x\perp}) = \mathbf{I},$
- $\mathbf{G}_{x\perp} \mathbf{X} = (\mathbf{A}_{x\perp}^\top \mathbf{A}_{x\perp})^{-1} \mathbf{A}^\top (\mathbf{I} - \mathbf{\Pi}_x) \mathbf{X} = \mathbf{0}.$

**Lemma E.5.** Assume  $m = 1$ , i.e.,  $\mathbf{A}$  is a column vector. We have

$$\mathbb{E}[\|\mathbf{K}\|^2] \leq 2d + 2\mathbb{E}\left[\frac{\|\mathbf{\Pi}_x \mathbf{A}\|^2 + \|\mu(\mathbf{X})\|^2}{\|(\mathbf{I} - \mathbf{\Pi}_x) \mathbf{A}\|^2}\right].$$

Similarly, it holds that

$$\|\mathbb{E}[\mathbf{L}]\|^2 \leq 2d + 2\mathbb{E}\left[\frac{\|\mathbf{\Pi}_x \mathbf{A}\|^2 + \|\mathbf{A}\|^2}{\|(\mathbf{I} - \mathbf{\Pi}_x) \mathbf{A}\|^2}\right].$$

*Proof.* We start by splitting the term into

$$\|\mathbf{K}\|^2 \leq 2\|\mathbf{X}\mathbf{H}_{a\perp}\|^2 + 2\|\mu(\mathbf{X})\mathbf{G}_{x\perp}\|^2.$$

Notice that  $\mathbf{H}_{a\perp} \mathbf{H}_{a\perp}^\top = (\mathbf{X}_{a\perp}^\top \mathbf{X}_{a\perp})^{-1}$  and similarly  $\mathbf{G}_{x\perp} \mathbf{G}_{x\perp}^\top = (\mathbf{A}_{x\perp}^\top \mathbf{A}_{x\perp})^{-1}$ .

Since  $\|\mathbf{M}\|^2 = \text{tr}(\mathbf{M}^\top \mathbf{M})$ , we have

$$\|\mathbf{K}\|^2 \leq 2\text{tr}((\mathbf{X}^\top \mathbf{X})(\mathbf{X}_{a\perp}^\top \mathbf{X}_{a\perp})^{-1}) + 2\text{tr}(\mu(\mathbf{X})^\top \mu(\mathbf{X})(\mathbf{A}_{x\perp}^\top \mathbf{A}_{x\perp})^{-1}).$$

By definition of  $\mathbf{A}_{x\perp}$ , when  $m = 1$ , we have

$$(\mathbf{A}_{x\perp}^\top \mathbf{A}_{x\perp})^{-1} = \frac{1}{\|(\mathbf{I} - \Pi_x)\mathbf{A}\|^2}.$$

For the term  $(\mathbf{X}_{a\perp}^\top \mathbf{X}_{a\perp})^{-1}$ , the Sherman–Morrison formula leads to

$$(\mathbf{X}_{a\perp}^\top \mathbf{X}_{a\perp})^{-1} = (\mathbf{X}^\top \mathbf{X})^{-1} + \frac{1}{1 - \mathbf{b}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{b}} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{b} \mathbf{b}^\top (\mathbf{X}^\top \mathbf{X})^{-1}$$

with  $\mathbf{b} = 1/\|\mathbf{A}\| \cdot \mathbf{X}^\top \mathbf{A} \in \mathbb{R}^d$ . Further simplifying, we obtain

$$\text{tr}((\mathbf{X}^\top \mathbf{X})(\mathbf{X}_{a\perp}^\top \mathbf{X}_{a\perp})^{-1}) = \text{tr}\left(\mathbf{I} + \frac{\Pi_x \mathbf{A} \mathbf{A}^\top \Pi_x}{\|\mathbf{A}\|^2 - \|\Pi_x \mathbf{A}\|^2}\right) = d + \frac{\|\Pi_x \mathbf{A}\|^2}{\|(\mathbf{I} - \Pi_x)\mathbf{A}\|^2}.$$

For the second part of the proof, we start by applying Jensen inequality:

$$\|\mathbb{E}[\mathbf{L}]\|^2 \leq \mathbb{E}[\|\mathbf{L}\|^2].$$

The rest of the proof follows along the same arguments, replacing  $\mu(\mathbf{X})$  by  $\mathbf{A}$ .  $\square$

## F RELATED WORK TABLE

Table 5: Comparison to related work.

METHOD	$p(\mathbf{a} \mathbf{x})$ REQUIRED	TRAINING	TEST COST	WEIGHT SHARING	APPROXIMATE $p(\mathbf{y} \mathbf{x})$
IMPUTATION	×	1 MODEL, 1 STEP	= NO PI	✓	×
DISTILLATION (LOPEZ-PAZ ET AL., 2015)	×	2 MODELS, 2 STEPS	= NO PI	×	×
HET. DROPOUT (LAMBERT ET AL., 2018)	×	1 MODEL, 1 STEP	= NO PI	✓	✓
MIML-FCN+ (YANG ET AL., 2017)	×	1 MODEL, 1 STEP	= NO PI	×	×
FULL MARGINALIZATION	✓	1 MODEL, 1 STEP	$\mathcal{O}(S * \text{NO PI})$	✓	✓
TRAM (OURS)	×	1 MODEL, 1 STEP	= NO PI	✓	✓
HET-TRAM (OURS)	×	1 MODEL, 1 STEP	= NO PI	✓	✓
DISTILLED-TRAM (OURS)	×	2 MODELS, 2 STEPS	= NO PI	✓	✓

## G TWO-STEP TRAM: IMAGENET SCALE REPRESENTATION LEARNING EXPERIMENT

We conduct experiment to test two things: 1) does the one-step TRAM procedure, introduced in §3.2 which is easier for practitioners to implement, approximate the two-step TRAM procedure well and 2) can the results of the toy representation learning experiment, §2.2 be replicated in a larger scale setting.

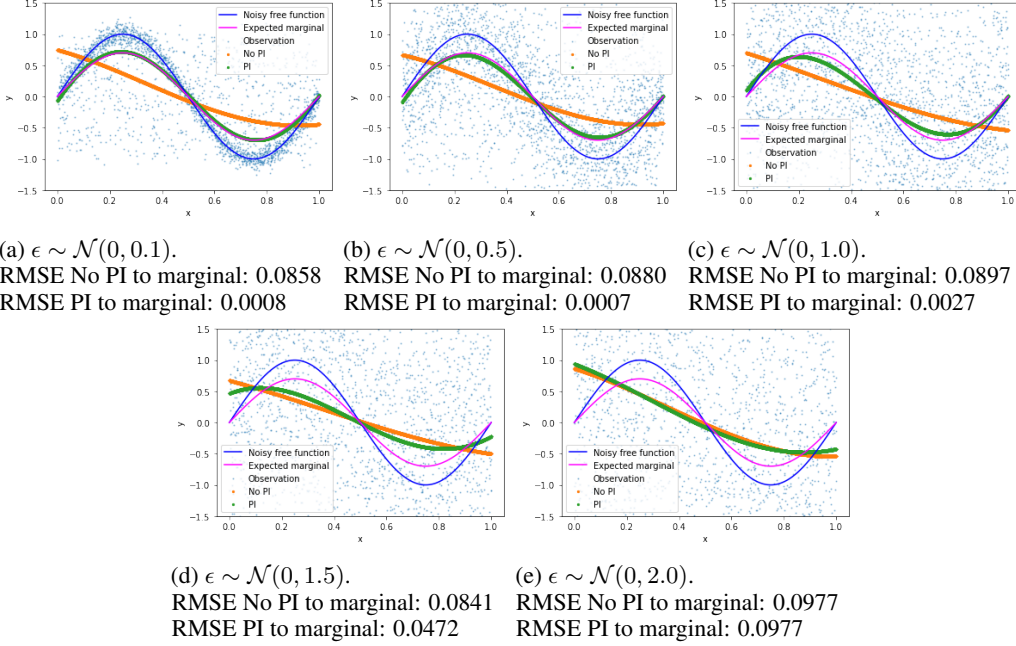
We train a feature extractor with and without access to PI on ImageNet, following the same procedure, architecture and dataset used in the main paper. We then freeze the feature extractor and train a single dense/linear layer with softmax activation on top of the fixed features. We then evaluate the efficacy of these features trained with and without PI using this “linear probe” evaluation widely used in the representation learning literature (Chen et al., 2020).

The results are presented in Table 6. We see that the simpler single-step TRAM method approximates the more complicated two-step TRAM method very well. In addition we see that the features learned by the network with access to PI which are then frozen and evaluated using a linear probe protocol perform better in terms of accuracy and log-likelihood.

## H TOY EXPERIMENT: VARY $\epsilon$

Table 6: Two-step TRAM: scaling up our toy representation learning experiment. ImageNet validation set negative log-likelihood and accuracy. Averaged over 10 training runs  $\pm 1$  std. dev.

METHOD	$\downarrow$ NLL	$\uparrow$ ACCURACY
ONE-STEP No PI	$1.264 \pm 0.007$	$71.7 \pm 0.2$
TWO-STEP No PI	$1.265 \pm 0.008$	$71.7 \pm 0.3$
ONE-STEP TRAM	$1.225 \pm 0.006$	$72.5 \pm 0.2$
TWO-STEP TRAM	$1.226 \pm 0.002$	$72.7 \pm 0.2$

Figure 4: Varying the influence of  $\epsilon$  on our motivating toy experiment.

We vary the standard deviation of  $\epsilon$  used in our motivating toy experiment, §2.2. The results can be seen graphically in Fig. 4. Fig. 4 also contains the average RMSE to the true marginal across the data points plotted. The graphical and numerical results demonstrate that even for large levels of noise PI aids with representation learning but as expected, as the level of noise grows the advantage of using PI diminishes as it becomes increasingly difficult to distinguish irreducible noise from noise which can be explained away with PI.

## I IMAGENET EXPERIMENT PI ABLATION

We run an ablation, removing PI feature: the probability of the label assigned by the model from the PI set. We are thus left with just one PI feature, the one-hot encoded ID of the model that produced the label.

We see the results in Table 7. As expected (and predicted by our theoretical analysis), removing informative PI reduces the effectiveness of TRAM. Nonetheless, TRAM with the reduced PI feature set still outperforms the No PI baseline, with accuracy and log-likelihood lying between the No PI and full PI feature set TRAM methods.



Table 7: ImageNet ablation with reduced PI feature set. ImageNet validation set negative log-likelihood and accuracy. Averaged over 10 training runs  $\pm$  1 std. dev.

METHOD	$\downarrow$ NLL	$\uparrow$ ACCURACY
NO PI	$1.264 \pm 0.007$	$71.7 \pm 0.2$
TRAM WITH FULL PI SET	$1.225 \pm 0.006$	$72.5 \pm 0.2$
TRAM WITH REDUCED PI SET	$1.246 \pm 0.004$	$72.0 \pm 0.2$