

A Related work

Self-supervised learning (SSL) methods in practice: There has been a flurry of self-supervised methods lately. One class of methods reconstruct images from corrupted or incomplete versions of it, like denoising auto-encoders [51], image inpainting [40], and split-brain autoencoder [57]. Pretext tasks are also created using visual common sense, including predicting rotation angle [18], relative patch position [12], recovering color channels [56], solving jigsaw puzzle games [38], and discriminating images created from distortion [13]. We refer to the above procedures as reconstruction-based SSL. Another popular paradigm is contrastive learning [9, 10]. The idea is to learn representations that bring similar data points closer while pushing randomly selected points further away [53, 33, 3] or to maximize a contrastive-based mutual information lower bound between different views [25, 39, 46]. A popular approach for text domain is based on language modeling where models like BERT and GPT create auxiliary tasks for next word predictions [11, 41]. The natural ordering or topology of data is also exploited in video-based [54, 37, 15], graph-based [55, 27] or map-based [58] self-supervised learning. For instance, the pretext task is to determine the correct temporal order for video frames as in [37].

Theory for self-supervised learning: Our work initiates some theoretical understanding on the reconstruction-based SSL. Related to our work is the recent theoretical analysis of contrastive learning. [3] shows guarantees for representations from contrastive learning on *linear classification* tasks using a class conditional independence assumption, but do not handle approximate conditional independence. Recently, (author?) [47] show that contrastive learning representations can *linearly* recover any continuous functions of the underlying topic posterior under a topic modeling assumption for text. While their assumption bears some similarity to ours, the assumption of independent sampling of words that they exploit is strong and not generalizable to other domains like images. More recently, concurrent work by [48] shows guarantees for contrastive learning, but not reconstruction-based SSL, with a multi-view redundancy assumptions that is very similar to our CI assumption. [52] theoretically studies contrastive learning on the hypersphere through intuitive properties like alignment and uniformity of representations; however there is no theoretical connection made to downstream tasks. There is a mutual information maximization view of contrastive learning, but [49] points out issues with it. Previous attempts to explain negative sampling [36] based methods use the theory of noise contrastive estimation [22, 34]. However, guarantees are only asymptotic and not for downstream tasks. CI is also used in sufficient dimension reduction [17, 16]. CI and redundancy assumptions on multiple views [31, 2] are used to analyze a canonical-correlation based dimension reduction algorithm. Finally, [1, 50] provide a theoretical analysis for denoising auto-encoder.

B Omitted Results with Conditional Independence

B.1 Warm-up: jointly Gaussian variables

We assume X_1, X_2, Y are jointly Gaussian, and so the optimal regression functions are all linear, i.e., $\mathbb{E}[Y|X_1] = \mathbb{E}^L[Y|X_1]$. We also assume data is centered: $\mathbb{E}[X_i] = 0$ and $\mathbb{E}[Y] = 0$. Non-centered data can easily be handled by learning an intercept. All relationships between random variables can then be captured by the (partial) covariance matrix. Therefore it is easy to quantify the CI property and establish the necessary and sufficient conditions that make X_2 a reasonable pretext task.

Assumption B.1. (Jointly Gaussian) X_1, X_2, Y are jointly Gaussian.

Assumption B.2. (Conditional independence) $X_1 \perp X_2 | Y$.

Claim B.1 (Closed-form solution). Under Assumption B.1, the representation function and optimal prediction that minimize the population risk can be expressed as follows:

$$\psi^*(\mathbf{x}_1) := \mathbb{E}^L[X_2|X_1 = \mathbf{x}_1] = \Sigma_{X_2 X_1} \Sigma_{X_1 X_1}^{-1} \mathbf{x}_1 \quad (1)$$

$$\text{Our target } f^*(\mathbf{x}_1) := \mathbb{E}^L[Y|X_1 = \mathbf{x}_1] = \Sigma_{Y X_1} \Sigma_{X_1 X_1}^{-1} \mathbf{x}_1. \quad (2)$$

Our prediction for downstream task with representation ψ^* will be: $g(\cdot) := \mathbb{E}^L[Y|\psi^*(X_1)]$. Recall from Equation 2 that the partial covariance matrix between X_1 and X_2 given Y is $\Sigma_{X_1 X_2|Y} \equiv \Sigma_{X_1 X_2} - \Sigma_{X_1 Y} \Sigma_{Y Y}^{-1} \Sigma_{Y X_2}$. This partial covariance matrix captures the correlation between X_1 and

359 X_2 given Y . For jointly Gaussian random variables, CI is equivalent to $\Sigma_{X_1 X_2|Y} = 0$. We first
 360 analyze the approximation error based on the property of this partial covariance matrix.

361 **Lemma B.2** (Approximation error). *Under Assumption B.1, B.2, if $\Sigma_{X_2 Y}$ has rank k , $e_{\text{apx}}(\psi^*) = 0$.*

362 **Remark B.1.** $\Sigma_{X_2 Y}$ being full column rank implies that $\mathbb{E}[X_2|Y]$ has rank k , i.e., X_2 depends on all
 363 directions of Y and thus captures all directions of information of Y . This is a necessary assumption
 364 for X_2 to be a reasonable pretext task for predicting Y . $e_{\text{apx}}(\psi^*) = 0$ means f^* is linear in ψ^* .
 365 Therefore ψ^* selects d_2 out of d_1 features that are sufficient to predict Y .

366 Next we consider the estimation error that characterizes the number of samples needed to learn a
 367 prediction function $f(\mathbf{x}_1) = \hat{\mathbf{W}}\psi^*(\mathbf{x}_1)$ that generalizes.

368 **Theorem B.3** (Estimation error). *Fix a failure probability $\delta \in (0, 1)$. Under Assumption B.1, B.2, if
 369 $n_2 \gg k + \log(1/\delta)$, excess risk of the learned predictor $\mathbf{x}_1 \rightarrow \hat{\mathbf{W}}\psi^*(\mathbf{x}_1)$ on the target task satisfies*

$$\text{ER}_{\psi^*}(\hat{\mathbf{W}}) \leq \mathcal{O}\left(\frac{\text{Tr}(\Sigma_{Y Y|X_1})(k + \log(k/\delta))}{n_2}\right),$$

370 *with probability at least $1 - \delta$.*

371 Here $\Sigma_{Y Y|X_1} \equiv \Sigma_{Y Y} - \Sigma_{Y X_1} \Sigma_{X_1 X_1}^{-1} \Sigma_{X_1 Y}$ captures the noise level and is the covariance matrix
 372 of tespeche residual term $Y - f^*(X_1) = Y - \Sigma_{Y X_1} \Sigma_{X_1 X_1}^{-1} X_1$. Compared to directly using X_1
 373 to predict Y , self-supervised learning reduces the sample complexity from $\tilde{\mathcal{O}}(d_1)$ to $\tilde{\mathcal{O}}(k)$. We
 374 generalize these results even when only a weaker form of CI holds.

375 **Assumption B.3** (Conditional independence given latent variables). *There exists some latent variable
 376 $Z \in \mathbb{R}^m$ such that $X_1 \perp X_2 | \bar{Y}$, and $\Sigma_{X_2 \bar{Y}}$ is of rank $k + m$, where $\bar{Y} = [Y, Z]$.*

377 This assumption lets introduce some reasonable latent variables that capture the information between
 378 X_1 and X_2 apart from Y . $\Sigma_{X_2 \bar{Y}}$ being full rank says that all directions of \bar{Y} are needed to predict
 379 X_2 , and therefore Z is not redundant. For instance, when $Z = X_1$, the assumption is trivially true
 380 but Z is not the minimal latent information we want to add. Note it implicitly requires $d_2 \geq k + m$.

381 **Corollary B.4.** *Under Assumption B.1, B.3, the approximation error $e_{\text{apx}}(\psi^*)$ is 0.*

382 Under CI with latent variable, we can generalize Theorem B.3 by replacing k by $k + m$.

383 C Some Useful Facts

384 C.1 Relation of Inverse Covariance Matrix and Partial Correlation

385 en For a covariance matrix of joint distribution for variables X, Y , the covariance matrix is

$$\begin{bmatrix} \Sigma_{XX} & \Sigma_{XY} \\ \Sigma_{YX} & \Sigma_{YY} \end{bmatrix} = \begin{bmatrix} \Sigma_{X_1 X_1} & \Sigma_{X_1 X_2} & \Sigma_{X_1 Y} \\ \Sigma_{X_2 X_1} & \Sigma_{X_2 X_2} & \Sigma_{X_2 Y} \\ \Sigma_{Y X_1} & \Sigma_{X_2 Y} & \Sigma_{Y Y} \end{bmatrix}.$$

386 Its inverse matrix Σ^{-1} satisfies

$$\Sigma^{-1} = \begin{bmatrix} \mathbf{A} & \rho \\ \rho^\top & B \end{bmatrix}.$$

387 Here $\mathbf{A}^{-1} = \Sigma_{XX} - \Sigma_{XY} \Sigma_{YY}^{-1} \Sigma_{YX} \equiv \text{cov}(X - \mathbb{E}^L[X|Y], X - \mathbb{E}^L[X|Y]) := \Sigma_{X X \cdot Y}$, the
 388 partial covariance matrix of X given Y .

389 C.2 Relation to Conditional Independence

390 *Proof of Lemma F.4.*

391 **Fact C.1.** *When $X_1 \perp X_2 | Y$, the partial covariance between X_1, X_2 given Y is 0:*

$$\begin{aligned} \Sigma_{X_1 X_2 \cdot Y} &:= \text{cov}(X_1 - \mathbb{E}^L[X_1|Y], X_2 - \mathbb{E}^L[X_2|Y]) \\ &\equiv \Sigma_{X_1 X_2} - \Sigma_{X_1 Y} \Sigma_{Y Y}^{-1} \Sigma_{Y X_2} = 0. \end{aligned}$$

392 The derivation comes from the following:

393 **Lemma C.1** (Conditional independence (Adapted from [28])). *For random variables X_1, X_2 and a*
 394 *random variable Y with finite values, conditional independence $X_1 \perp X_2 | Y$ is equivalent to:*

$$\sup_{f \in N_1, g \in N_2} \mathbb{E}[f(X_1)g(X_2)|Y] = 0. \quad (3)$$

395 Here $N_i = \{f : \mathbb{R}^{d_i} \rightarrow \mathbb{R} : \mathbb{E}[f(X_i)|Y] = 0\}$, $i = 1, 2$.

396 Notice for arbitrary function f , $\mathbb{E}[f(X)|Y] = \mathbb{E}^L[f(X)|\phi_y(Y)]$ with one-hot encoding of discrete
 397 variable Y . Therefore for any feature map we can also get that conditional independence ensures:

$$\begin{aligned} \Sigma_{\phi_1(X_1)\phi_2(X_2)|Y} &:= \text{cov}(\phi_1(X_1) - \mathbb{E}^L[\phi_1(X_1)|\phi_y(Y)], \phi_2(X_2) - \mathbb{E}^L[\phi_2(X_2)|\phi_y(Y)]) \\ &= \mathbb{E}[\bar{\phi}_1(X_1)\bar{\phi}_2(X_2)^\top] = 0. \end{aligned}$$

398 Here $\bar{\phi}_1(X_1) = \phi_1(X_1) - \mathbb{E}[\phi_1(X_1)|\phi_y(Y)]$ is mean zero given Y , and vice versa for $\bar{\phi}_2(X_2)$. This
 399 thus finishes the proof for Lemma F.4. \square

400 C.3 Technical Facts for Matrix Concentration

401 We include this covariance concentration result that is adapted from Claim A.2 in [14]:

402 **Claim C.2** (covariance concentration for gaussian variables). *Let $\mathbf{X} = [x_1, x_2, \dots, x_n]^\top \in \mathbb{R}^{n \times d}$*
 403 *where each $x_i \sim \mathcal{N}(0, \Sigma_X)$. Suppose $n \gg k + \log(1/\delta)$ for $\delta \in (0, 1)$. Then for any given matrix*
 404 *$B \in \mathbb{R}^{d \times m}$ that is of rank k and is independent of \mathbf{X} , with probability at least $1 - \frac{\delta}{10}$ over \mathbf{X} we*
 405 *have*

$$0.9B^\top \Sigma_X B \preceq \frac{1}{n} B^\top \mathbf{X}^\top \mathbf{X} B \preceq 1.1B^\top \Sigma_X B. \quad (4)$$

406 And we will also use Claim A.2 from [14] for concentrating subgaussian random variable.

407 **Claim C.3** (covariance concentration for subgaussian variables). *Let $\mathbf{X} = [x_1, x_2, \dots, x_n]^\top \in \mathbb{R}^{n \times d}$*
 408 *where each $x_i \sim \mathcal{N}(0, \Sigma_X)$. Suppose $n \gg \rho^4(k + \log(1/\delta))$ for $\delta \in (0, 1)$. Then for any given*
 409 *matrix $B \in \mathbb{R}^{d \times m}$ that is of rank k and is independent of \mathbf{X} , with probability at least $1 - \frac{\delta}{10}$ over \mathbf{X}*
 410 *we have*

$$0.9B^\top \Sigma_X B \preceq \frac{1}{n} B^\top \mathbf{X}^\top \mathbf{X} B \preceq 1.1B^\top \Sigma_X B. \quad (5)$$

411 **Claim C.4.** *Let $\mathbf{Z} \in \mathbb{R}^{n \times k}$ be a matrix with row vectors sampled from i.i.d Gaussian distribution*
 412 *$\mathcal{N}(0, \Sigma_Z)$. Let $P \in \mathbb{R}^{n \times n}$ be a fixed projection onto a space of dimension d . Then with a fixed*
 413 *$\delta \in (0, 1)$, we have:*

$$\|P\mathbf{Z}\|_F^2 \lesssim \text{Tr}(\Sigma_Z)(d + \log(k/\delta)),$$

414 *with probability at least $1 - \delta$.*

415 *Proof of Claim C.4.* Each t -th column of \mathbf{Z} is an n -dim vector that is i.i.d sampled from Gaussian
 416 distribution $\mathcal{N}(0, \Sigma_{tt})$.

$$\begin{aligned} \|P\mathbf{Z}\|_F^2 &= \sum_{t=1}^k \|P\mathbf{z}_t\|^2 \\ &= \sum_{t=1}^k \mathbf{z}_t^\top P \mathbf{z}_t. \end{aligned}$$

417 Each term satisfy $\Sigma_{kk}^{-1} \|P\mathbf{z}_t\|^2 \sim \chi^2(d)$, and therefore with probability at least $1 - \delta'$ over \mathbf{z}_t ,

$$\Sigma_{kk}^{-1} \|P\mathbf{z}_t\|^2 \lesssim d + \log(1/\delta').$$

418 Using union bound, take $\delta' = \delta/k$ and summing over $t \in [k]$ we get:

$$\|P\mathbf{Z}\|_F^2 \lesssim \text{Tr}(\Sigma_Z)(d + \log(k/\delta)).$$

419 \square

420 **Theorem C.5** (Hanson-Wright Inequality (Theorem 1.1 from [43])). Let $X = (X_1, X_2, \dots, X_n) \in$
421 \mathbb{R}^n be a random vector with independent components X_i which satisfy $\mathbb{E}[X_i] = 0$ and $\|X_i\|_{\psi_2} \leq K$.
422 Let A be an $n \times n$ matrix. Then, for every $t \geq 0$,

$$\mathbb{P}[|X^\top A X - \mathbb{E}[X^\top A X]| > t] \leq 2 \exp \left\{ -c \min \left(\frac{t^2}{K^4 \|A\|_F^2}, \frac{t}{K^2 \|A\|} \right) \right\}.$$

Theorem C.6 (Vector Bernstein Inequality (Theorem 12 in [21])). Let X_1, \dots, X_m be independent zero-mean vector-valued random variables. Let

$$N = \left\| \sum_{i=1}^m X_i \right\|_2.$$

423 Then

$$\mathbb{P}[N \geq \sqrt{V} + t] \leq \exp \left(\frac{-t^2}{4V} \right),$$

424 where $V = \sum_i \mathbb{E} \|X_i\|_2^2$ and $t \leq V/(\max \|X_i\|_2)$.

425 **Lemma C.7.** Let $\mathbf{Z} \in \mathbb{R}^{n \times k}$ be a matrix whose row vectors are n independent mean-zero (condi-
426 tional on P) σ -sub-Gaussian random vectors. With probability $1 - \delta$:

$$\|P\mathbf{Z}\|^2 \lesssim \sigma^2(d + \log(d/\delta)).$$

427 *Proof of Lemma C.7.* Write $P = \mathbf{U}\mathbf{U}^\top = [\mathbf{u}_1, \dots, \mathbf{u}_d]$ where \mathbf{U} is orthogonal matrix in $\mathbb{R}^{n \times d}$
428 where $\mathbf{U}^\top \mathbf{U} = I$.

$$\begin{aligned} \|P\mathbf{Z}\|_F^2 &= \|\mathbf{U}^\top \mathbf{Z}\|_F^2 \\ &= \sum_{j=1}^d \|\mathbf{u}_j^\top \mathbf{Z}\|^2 \\ &= \sum_{j=1}^d \left\| \sum_{i=1}^n \mathbf{u}_{ji} \mathbf{z}_i \right\|^2, \end{aligned}$$

429 where each $\mathbf{z}_i \in \mathbb{R}^k$ being the i -th row of \mathbf{Z} is a centered independent σ sub-Gaussian random
430 vectors. To use vector Bernstein inequality, we let $X := \sum_{i=1}^n X_i$ with $X_i := \mathbf{u}_{ji} \mathbf{z}_i$. We have X_i is
431 zero mean: $\mathbb{E}[X_i] = \mathbb{E}[\mathbf{u}_{ji} \mathbb{E}[\mathbf{z}_i | \mathbf{u}_{ji}]] = \mathbb{E}[\mathbf{u}_{ji} \cdot 0] = 0$.

$$\begin{aligned} V &:= \sum_i \mathbb{E} \|\mathbf{z}_i\|_2^2 \\ &= \sum_i \mathbb{E}[\mathbf{u}_{ji}^2 \mathbf{z}_i^\top \mathbf{z}_i] \\ &= \sum_i \mathbb{E}_{\mathbf{u}_{ji}}[\mathbf{u}_{ji}^2 \mathbb{E}[\|\mathbf{z}_i\|_2^2 | \mathbf{u}_{ji}]] \\ &\leq \sigma^2 \sum_i \mathbb{E}_{\mathbf{u}_{ji}}[\mathbf{u}_{ji}^2] \\ &= \sigma^2. \end{aligned}$$

432 Therefore by vector Bernstein Inequality, with probability at least $1 - \delta/d$, $\|X\| \leq \sigma(1 + \sqrt{\log(d/\delta)})$.
433 Then by taking union bound, we get that $\|P\mathbf{Z}\|^2 = \sum_{j=1}^d \|\mathbf{u}_j^\top \mathbf{Z}\|^2 \lesssim \sigma^2(d + \log(d/\delta))$ with
434 probability $1 - \delta$.

435 □

436 **Corollary C.8.** Let $\mathbf{Z} \in \mathbb{R}^{n \times k}$ be a matrix whose row vectors are n independent samples from
437 centered (conditioned on P) multinomial probabilities (p_1, p_2, \dots, p_k) (where p_t could be different
438 across each row). Let $P \in \mathbb{R}^{n \times n}$ be a projection onto a space of dimension d (that might be
439 dependent with \mathbf{Z}). Then we have

$$\|P\mathbf{Z}\|^2 \lesssim d + \log(d/\delta).$$

440 with probability $1 - \delta$.

441 D Omitted Proofs with Conditional Independence

Proof of Lemma B.2.

$$\text{cov}(X_1|Y, X_2|Y) = \Sigma_{X_1X_2} - \Sigma_{X_1Y} \Sigma_{YY}^{-1} \Sigma_{YX_2} = 0.$$

442 By plugging it into the expression of $\mathbb{E}^L[X_2|X_1]$, we get that

$$\begin{aligned} \psi(x_1) &:= \mathbb{E}^L[X_2|X_1 = x_1] = \Sigma_{X_2X_1} \Sigma_{X_1X_1}^{-1} x_1 \\ &= \Sigma_{X_2Y} \Sigma_{YY}^{-1} \Sigma_{YX_1} \Sigma_{X_1X_1}^{-1} x_1 \\ &= \Sigma_{X_2Y} \Sigma_{YY}^{-1} \mathbb{E}^L[Y|X_1]. \end{aligned}$$

443 Therefore, as long as Σ_{X_2Y} of rank k , it has left inverse matrix and we get: $\mathbb{E}^L[Y|X_1 = x_1] =$
444 $\Sigma_{X_2Y}^\dagger \Sigma_{YY} \psi(x_1)$. Therefore there's no approximation error in using ψ to predict Y .

445 □

446 *Proof of Corollary B.4.* Let selector operator S_y be the mapping such that $S_y \bar{Y} = Y$, we overload
447 it as the matrix that ensure $S_y \Sigma_{\bar{Y}X} = \Sigma_{YX}$ for any random variable X as well.

448 From Lemma B.2 we get that there exists W such that $\mathbb{E}^L[\bar{Y}|X_1] = W \mathbb{E}^L[X_2|X_1]$, just plugging in
449 S_y we get that $\mathbb{E}^L[Y|X_1] = (S_y W) \mathbb{E}^L[X_2|X_1]$.

450 □

451 *Proof of Theorem B.3.* Since N is mean zero, $f^*(X_1) = \mathbb{E}[Y|X_1] = (\mathbf{A}^*)^\top X_1$.

452 $\mathbb{E}^L[Y|X_1 = x_1] = \Sigma_{X_2Y}^\dagger \Sigma_{YY} \psi(x_1)$. Let $\mathbf{W}^* = \Sigma_{YY} \Sigma_{YX_2}^\dagger$.

453 First we have the basic inequality,

$$\begin{aligned} \frac{1}{2n_2} \|\mathbf{Y} - \psi(\mathbf{X}_1) \hat{\mathbf{W}}\|_F^2 &\leq \frac{1}{2n_2} \|\mathbf{Y} - \mathbf{X}_1 \mathbf{A}^*\|_F^2 \\ &= \frac{1}{2n_2} \|\mathbf{Y} - \psi(\mathbf{X}_1) \mathbf{W}^*\|_F^2. \end{aligned}$$

454 Therefore

$$\begin{aligned} \|\psi(\mathbf{X}_1) \mathbf{W}^* - \psi(\mathbf{X}_1) \hat{\mathbf{W}}\|^2 &\leq 2 \langle N, \psi(\mathbf{X}_1) \mathbf{W}^* - \psi(\mathbf{X}_1) \hat{\mathbf{W}} \rangle \\ &= 2 \langle P_{\psi(\mathbf{X}_1)} \mathbf{N}, \psi(\mathbf{X}_1) \mathbf{W}^* - \psi(\mathbf{X}_1) \hat{\mathbf{W}} \rangle \\ &\leq 2 \|P_{\psi(\mathbf{X}_1)} \mathbf{N}\|_F \|\psi(\mathbf{X}_1) \mathbf{W}^* - \psi(\mathbf{X}_1) \hat{\mathbf{W}}\|_F \\ \Rightarrow \|\psi(\mathbf{X}_1) \mathbf{W}^* - \psi(\mathbf{X}_1) \hat{\mathbf{W}}\| &\leq 2 \|P_{\psi(\mathbf{X}_1)} \mathbf{N}\|_F \\ &\lesssim \sqrt{\text{Tr}(\Sigma_{YY|X_1})(k + \log k/\delta)}. \quad (\text{from Claim C.4}) \end{aligned}$$

The last inequality is derived from Claim C.7 and the fact that each row of \mathbf{N} follows gaussian distribution $\mathcal{N}(0, \Sigma_{YY|X_1})$. Therefore

$$\frac{1}{n_2} \|\psi(\mathbf{X}_1) \mathbf{W}^* - \psi(\mathbf{X}_1) \hat{\mathbf{W}}\|_F^2 \lesssim \frac{\text{Tr}(\Sigma_{YY|X_1})(k + \log k/\delta)}{n_2}.$$

455 Next we need to concentrate $1/n \mathbf{X}_1^\top \mathbf{X}_1$ to Σ_X . Suppose $\mathbb{E}^L[X_2|X_1] = \mathbf{B}^\top X_1$, i.e., $\phi(x_1) =$
456 $\mathbf{B}^\top x_1$, and $\phi(\mathbf{X}_1) = \mathbf{X}_1 \mathbf{B}$. With Claim C.2 we have $1/n \phi(\mathbf{X}_1)^\top \phi(\mathbf{X}_1) = 1/n \mathbf{B}^\top \mathbf{X}_1^\top \mathbf{X}_1 \mathbf{B}$
457 satisfies:

$$0.9 \mathbf{B}^\top \Sigma_X \mathbf{B} \preceq 1/n_2 \phi(\mathbf{X}_1)^\top \phi(\mathbf{X}_1) \preceq 1.1 \mathbf{B}^\top \Sigma_X \mathbf{B}$$

458 Therefore we also have:

$$\begin{aligned} &\mathbb{E}[(\mathbf{W}^* - \hat{\mathbf{W}})^\top \psi(x_1)] \\ &= \|\Sigma_X^{1/2} \mathbf{B}(\mathbf{W}^* - \hat{\mathbf{W}})\|_F^2 \\ &\leq \frac{1}{0.9 n_2 k} \|\psi(\mathbf{X}_1) \mathbf{W}^* - \psi(\mathbf{X}_1) \hat{\mathbf{W}}\|_F^2 \lesssim \frac{\text{Tr}(\Sigma_{YY|X_1})(k + \log k/\delta)}{n_2}. \end{aligned}$$

459 □

460 D.1 Omitted Proof for General Random Variables

461 *Proof of Lemma 3.1.* Let the representation function ψ be defined as:

$$\begin{aligned}\psi(\cdot) &:= \mathbb{E}[X_2|X_1] = \mathbb{E}[\mathbb{E}[X_2|X_1, Y]|X_1] \\ &= \mathbb{E}[\mathbb{E}[X_2|Y]|X_1] \quad (\text{uses CI}) \\ &= \sum_y P(Y = y|X_1) \mathbb{E}[X_2|Y = y] \\ &=: f(X_1)^\top \mathbf{A},\end{aligned}$$

462 where $f : \mathbb{R}^{d_1} \rightarrow \Delta_{\mathcal{Y}}$ satisfies $f(x_1)_y = P(Y = y|X_1 = x_1)$, and $\mathbf{A} \in \mathbb{R}^{\mathcal{Y} \times d_2}$ satisfies $\mathbf{A}_{y,:} =$
463 $\mathbb{E}[X_2|Y = y]$. Here Δ_d denotes simplex of dimension d , which represents the discrete probability
464 density over support of size d .

465 Let $\mathbf{B} = \mathbf{A}^\dagger \in \mathbb{R}^{\mathcal{Y} \times d_2}$ be the pseudoinverse of matrix \mathbf{A} , and we get $\mathbf{B}\mathbf{A} = \mathbf{I}$ from our assumption
466 that \mathbf{A} is of rank $|\mathcal{Y}|$. Therefore $f(x_1) = \mathbf{B}\psi(x_1), \forall x_1$. Next we have:

$$\begin{aligned}\mathbb{E}[Y|X_1 = x_1] &= \sum_y P(Y = y|X_1 = x_1) \times y \\ &= \mathbf{Y}f(x_1) \\ &= (\mathbf{Y}\mathbf{B}) \cdot \psi(X_1).\end{aligned}$$

467 Here we denote by $\mathbf{Y} \in \mathbb{R}^{k \times \mathcal{Y}}$, $\mathbf{Y}_{:,y} = y$ that spans the whole support \mathcal{Y} . Therefore let $\mathbf{W}^* = \mathbf{Y}\mathbf{B}$
468 will finish the proof. □

470 *Proof of Theorem 3.2.* With Lemma 3.1 we know $e_{\text{apx}} = 0$, and therefore $\mathbf{W}^*\psi(X_1) \equiv f^*(X_1)$.
471 Next from basic inequality and the same proof as in Theorem B.3 we have:

$$\|\psi(\mathbf{X}_1)\mathbf{W}^* - \psi(\mathbf{X}_1)\hat{\mathbf{W}}\| \leq 2\|P_{\psi(\mathbf{X}_1)}\mathbf{N}\|_F$$

472 Notice \mathbf{N} is a random noise matrix whose row vectors are independent samples from some centered
473 distribution. Also we assumed $\mathbb{E}[\|\mathbf{N}\|^2|\mathbf{X}_1] \leq \sigma^2$, i.e. $\mathbb{E}[\|\mathbf{N}\|^2|\mathbf{N}] \leq \sigma^2$. Also, $P_{\psi(\mathbf{X}_1)}$ is a
474 projection to dimension c . From Lemma C.7 we have:

$$\|f^*(\mathbf{X}_1) - \psi(\mathbf{X}_1)\hat{\mathbf{W}}\| \leq \sigma\sqrt{c + \log c/\delta}.$$

475 Next, with Claim C.3 we have when $n \gg \rho^4(c + \log(1/\delta))$, since $\mathbf{W}^* - \hat{\mathbf{W}} \in \mathbb{R}^{d_2 \times k}$,

$$\begin{aligned}&0.9(\mathbf{W}^* - \hat{\mathbf{W}})^\top \Sigma_\psi(\mathbf{W}^* - \hat{\mathbf{W}}) \\ &\preceq \frac{1}{n_2}(\mathbf{W}^* - \hat{\mathbf{W}})^\top \sum_i \psi(x_1^{(i)})\psi(x_1^{(i)})^\top (\mathbf{W}^* - \hat{\mathbf{W}}) \preceq 1.1(\mathbf{W}^* - \hat{\mathbf{W}})^\top \Sigma_\psi(\mathbf{W}^* - \hat{\mathbf{W}})\end{aligned}$$

476 And therefore we could easily conclude that:

$$\mathbb{E}\|\hat{\mathbf{W}}^\top \psi(X_1) - f^*(X_1)\|^2 \lesssim \sigma^2 \frac{c + \log(c/\delta)}{n_2}.$$

477 □

478 D.2 Omitted proof of linear model with approximation error

479 *Proof of Theorem 3.5.* First we note that $Y = f^*(X_1) + N$, where $\mathbb{E}[N|X_1] = 0$ but $Y - (\mathbf{A}^*)^\top X_1$
480 is not necessarily mean zero, and this is where additional difficulty lies. Write approximation error
481 term $a(X_1) := f^*(X_1) - (\mathbf{A}^*)^\top X_1$, namely $Y = a(X_1) + (\mathbf{A}^*)^\top X_1 + N$. Also, $(\mathbf{A}^*)^\top X_1 \equiv$
482 $(\mathbf{W}^*)^\top \psi(X_1)$ with conditional independence.

483 Second, with KKT condition on the training data, we know that $\mathbb{E}[a(X_1)X_1^\top] = 0$.

484 Recall $\hat{\mathbf{W}} = \arg \min_{\mathbf{W}} \|\mathbf{Y} - \psi(\mathbf{X}_1)\mathbf{W}\|_F^2$. We have the basic inequality,

$$\begin{aligned} \frac{1}{2n_2} \|\mathbf{Y} - \psi(\mathbf{X}_1)\hat{\mathbf{W}}\|_F^2 &\leq \frac{1}{2n_2} \|\mathbf{Y} - \mathbf{X}_1\mathbf{A}^*\|_F^2 \\ &= \frac{1}{2n_2} \|\mathbf{Y} - \psi(\mathbf{X}_1)\mathbf{W}^*\|_F^2. \end{aligned}$$

$$\text{i.e., } \frac{1}{2n_2} \|\psi(\mathbf{X}_1)\mathbf{W}^* + \mathbf{a}(\mathbf{X}_1) + \mathbf{N} - \psi(\mathbf{X}_1)\hat{\mathbf{W}}\|_F^2 \leq \frac{1}{2n_2} \|\mathbf{a}(\mathbf{X}_1) + \mathbf{N}\|_F^2.$$

485 Therefore

$$\begin{aligned} &\frac{1}{2n_2} \|\psi(\mathbf{X}_1)\mathbf{W}^* - \psi(\mathbf{X}_1)\hat{\mathbf{W}}\|^2 \\ &\leq -\frac{1}{n_2} \langle \mathbf{a}(\mathbf{X}_1) + \mathbf{N}, \psi(\mathbf{X}_1)\mathbf{W}^* - \psi(\mathbf{X}_1)\hat{\mathbf{W}} \rangle \\ &= -\frac{1}{n_2} \langle \mathbf{a}(\mathbf{X}_1), \psi(\mathbf{X}_1)\mathbf{W}^* - \psi(\mathbf{X}_1)\hat{\mathbf{W}} \rangle - \langle \mathbf{N}, \psi(\mathbf{X}_1)\mathbf{W}^* - \psi(\mathbf{X}_1)\hat{\mathbf{W}} \rangle \end{aligned} \quad (6)$$

486 With Assumption 3.3 and by concentration $0.9 \frac{1}{n_2} \mathbf{X}_1 \mathbf{X}_1^\top \preceq \Sigma_{X_1} \preceq 1.1 \frac{1}{n_2} \mathbf{X}_1 \mathbf{X}_1^\top$, we have

$$\frac{1}{\sqrt{n_2}} \|\mathbf{a}(\mathbf{X}_1) \mathbf{X}_1^\top \Sigma_{X_1}^{-1/2}\|_F \leq 1.1 b_0 \sqrt{k} \quad (7)$$

487 Denote $\psi(\mathbf{X}_1) = \mathbf{X}_1 \mathbf{B}$, where $\mathbf{B} = \Sigma_{X_1}^{-1} \Sigma_{X_1 X_2}$ is rank k under exact CI since $\Sigma_{X_1 X_2} =$

488 $\Sigma_{X_1 Y} \Sigma_Y^{-1} \Sigma_{Y X_2}$. We have

$$\begin{aligned} &\frac{1}{n_2} \langle \mathbf{a}(\mathbf{X}_1), \psi(\mathbf{X}_1)\mathbf{W}^* - \psi(\mathbf{X}_1)\hat{\mathbf{W}} \rangle \\ &= \frac{1}{n_2} \langle \mathbf{a}(\mathbf{X}_1), \mathbf{X}_1 \mathbf{B} \mathbf{W}^* - \mathbf{X}_1 \mathbf{B} \hat{\mathbf{W}} \rangle \\ &= \frac{1}{n_2} \langle \Sigma_{X_1}^{-1/2} \mathbf{X}_1^\top \mathbf{a}(\mathbf{X}_1), \Sigma_{X_1}^{1/2} (\mathbf{B} \mathbf{W}^* - \mathbf{B} \hat{\mathbf{W}}) \rangle \\ &\leq \sqrt{\frac{k}{n_2}} \|\Sigma_{X_1}^{1/2} (\mathbf{B} \mathbf{W}^* - \mathbf{B} \hat{\mathbf{W}})\|_F \end{aligned} \quad (\text{from Ineq. (7)})$$

489 Back to Eqn. (6), we get

$$\begin{aligned} &\frac{1}{2n_2} \|\psi(\mathbf{X}_1)\mathbf{W}^* - \psi(\mathbf{X}_1)\hat{\mathbf{W}}\|_F^2 \\ &\lesssim \sqrt{\frac{k}{n_2}} \|\Sigma_{X_1}^{1/2} (\mathbf{B} \mathbf{W}^* - \mathbf{B} \hat{\mathbf{W}})\|_F + \frac{1}{n_2} \|P_{\mathbf{X}_1} \mathbf{N}\|_F \|\mathbf{X}_1 (\mathbf{B} \mathbf{W}^* - \mathbf{B} \hat{\mathbf{W}})\|_F \\ &\lesssim \left(\frac{\sqrt{k}}{n_2} + \frac{1}{n_2} \|P_{\mathbf{X}_1} \mathbf{N}\|_F \right) \|\mathbf{X}_1 (\mathbf{B} \mathbf{W}^* - \mathbf{B} \hat{\mathbf{W}})\|_F \\ &\implies \frac{1}{\sqrt{n_2}} \|\psi(\mathbf{X}_1)\mathbf{W}^* - \psi(\mathbf{X}_1)\hat{\mathbf{W}}\|_F \lesssim \sqrt{\frac{k + \log k/\delta}{n_2}}. \end{aligned}$$

490 Finally, by concentration we transfer the result from empirical loss to excess risk and get:

$$\mathbb{E}[\|\psi(\mathbf{X}_1)\mathbf{W}^* - \psi(\mathbf{X}_1)\hat{\mathbf{W}}\|^2] \lesssim \frac{k + \log(k/\delta)}{n_2}.$$

491

□

492 D.3 Argument on Denoising Auto-encoder or Context Encoder

493 **Remark D.1.** We note that since $X_1 \perp X_2 | Y$ ensures $X_1 \perp h(X_2) | Y$ for any deterministic function h ,
 494 we could replace X_2 by $h(X_2)$ and all results hold. Therefore in practice, we could use $h(\psi(X_1))$
 495 instead of $\psi(X_1)$ for downstream task. Specifically with denoising auto-encoder or context encoder,
 496 one could think about h as the inverse of decoder D ($h = D^{-1}$) and use $D^{-1}\psi \equiv E$ the encoder
 497 function as the representation for downstream tasks, which is more commonly used in practice.

498 This section explains what we claim in Remark D.1. For context encoder, the reconstruction loss
 499 targets to find the encoder E^* and decoder D^* that achieve

$$\min_E \min_D \mathbb{E} \|X_2 - D(E(X_1))\|_F^2, \quad (8)$$

500 where X_2 is the masked part we want to recover and X_1 is the remainder.

501 If we naively apply our theorem we should use $D^*(E^*(\cdot))$ as the representation, while in practice we
 502 instead use only the encoder part $E^*(\cdot)$ as the learned representation. We argue that our theory also
 503 support this practical usage if we view the problem differently. Consider the pretext task to predict
 504 $(D^*)^{-1}(X_2)$ instead of X_2 directly, namely,

$$\bar{E} \leftarrow \arg \min_E \mathbb{E} \|(D^*)^{-1}(X_2) - E(X_1)\|^2, \quad (9)$$

and then we should indeed use $E(X_1)$ as the representation. On one hand, when $X_1 \perp X_2 | Y$, it
 also satisfies $X_1 \perp (D^*)^{-1}(X_2) | Y$ since $(D^*)^{-1}$ is a deterministic function of X_2 and all our theory
 applies. On the other hand, the optimization on (8) or (9) give us similar result. Let

$$E^* = \arg \min_E \mathbb{E} [\|X_2 - D^*(E(X_1))\|^2],$$

505 and $\mathbb{E} \|X_2 - D^*(E^*(X_1))\|^2 \leq \epsilon$, then with pretext task as in (9) we have that:

$$\begin{aligned} \mathbb{E} \|(D^*)^{-1}(X_2) - E^*(X_1)\|^2 &= \mathbb{E} \|(D^*)^{-1}(X_2) - (D^*)^{-1} \circ D^*(E^*(X_1))\|^2 \\ &\leq \|(D^*)^{-1}\|_{\text{Lip}}^2 \mathbb{E} \|X_2 - D^*(E^*(X_1))\|^2 \\ &\leq L^2 \epsilon, \end{aligned}$$

506 where $L := \|(D^*)^{-1}\|_{\text{Lip}}$ is the Lipschitz constant for function $(D^*)^{-1}$. This is to say, in practice,
 507 we optimize over (8), and achieves a good representation $E^*(X_1)$ such that $\epsilon_{\text{pre}} \leq L\sqrt{\epsilon}$ and thus
 508 performs well for downstream tasks. (Recall ϵ_{pre} is defined in Theorem E.3 that measures how well
 509 we have learned the pretext task.)

510 E Beyond conditional independence

511 In the previous section, we focused on the case where exact CI is satisfied. A weaker but more
 512 practical assumption is that Y captures some portion of the dependence between X_1 and X_2 but not
 513 all. We start with the jointly-Gaussian case, where approximate CI is quantified by partial covariance
 514 matrix. We then generalize the results and introduce covariance operator to measure approximate CI.

515 E.1 Warm-up: Jointly Gaussian Variables

516 As before, for simplicity we assume all data is centered in this case.

Assumption E.1 (Approximate Conditional Independent Given Latent Variables). *Assume there exists some latent variable $Z \in \mathbb{R}^m$ such that*

$$\|\Sigma_{X_1}^{-1/2} \Sigma_{X_1, X_2 | \bar{Y}}\|_F \leq \epsilon_{CI},$$

517 $\sigma_{k+m}(\Sigma_{Y\bar{Y}}^\dagger \Sigma_{\bar{Y}X_2}) = \beta > 0$ ¹ and $\Sigma_{X_2, \bar{Y}}$ is of rank $k + m$, where $\bar{Y} = [Y, Z]$.

518 When X_1 is not exactly CI of X_2 given Y and Z , the approximation error depends on the norm of
 519 $\|\Sigma_{X_1}^{-1/2} \Sigma_{X_1, X_2 | \bar{Y}}\|_2$. Let \hat{W} be the solution from Equation ??.

520 **Theorem E.1.** *Under Assumption E.1 with constant ϵ_{CI} and β , then the excess risk satisfies*

$$\text{ER}_{\psi^*}[\hat{W}] := \mathbb{E}[\|\hat{W}^\top \psi^*(X_1) - f^*(X_1)\|_F^2] \lesssim \frac{\epsilon_{CI}^2}{\beta^2} + \text{Tr}(\Sigma_{Y|X_1}) \frac{d_2 + \log(d_2/\delta)}{n_2}.$$

521 *Proof of Theorem E.1.* Let $V := f^*(X_1) \equiv X_1 \Sigma_{X_1 X_1}^{-1} \Sigma_{X_1 Y}$ be our target direction. Denote the
 522 optimal representation matrix by $\Psi := \psi(X_1) \equiv X_1 A$ (where $A := \Sigma_{X_1 X_1}^{-1} \Sigma_{X_1 X_2}$).

¹ $\sigma_k(A)$ denotes k -th singular value of A , and A^\dagger is the pseudo-inverse of A .

Next we will make use of the conditional covariance matrix:

$$\Sigma_{X_1 X_2 | \bar{Y}} := \Sigma_{X_1 X_2} - \Sigma_{X_1 \bar{Y}} \Sigma_{\bar{Y}}^{-1} \Sigma_{\bar{Y} X_2},$$

and plug it in into the definition of Ψ :

$$\begin{aligned} \Psi &= \mathbf{X}_1 \Sigma_{X_1 X_1}^{-1} \Sigma_{X_1 \bar{Y}} \Sigma_{\bar{Y}}^{-1} \Sigma_{\bar{Y} X_2} + \mathbf{X}_1 \Sigma_{X_1 X_1}^{-1} \Sigma_{X_1 X_2 | \bar{Y}} \\ &=: \mathbf{L} + \mathbf{E}, \end{aligned}$$

where $\mathbf{L} := \mathbf{X}_1 \Sigma_{X_1 X_1}^{-1} \Sigma_{X_1 \bar{Y}} \Sigma_{\bar{Y}}^{-1} \Sigma_{\bar{Y} X_2}$ and $\mathbf{E} := \mathbf{X}_1 \Sigma_{X_1 X_1}^{-1} \Sigma_{X_1 X_2 | \bar{Y}}$. We analyze these two terms respectively.

For \mathbf{L} , we note that $\text{span}(\mathbf{V}) \subseteq \text{span}(\mathbf{L})$: $\mathbf{L} \Sigma_{X_2 \bar{Y}}^\dagger \Sigma_{\bar{Y}} = \mathbf{X}_1 \Sigma_{X_1 X_1}^{-1} \Sigma_{X_1 \bar{Y}}$. By right multiplying the selector matrix S_Y we have: $\mathbf{L} \Sigma_{X_2 \bar{Y}}^\dagger \Sigma_{\bar{Y} Y} = \mathbf{X}_1 \Sigma_{X_1 X_1}^{-1} \Sigma_{X_1 Y}$, i.e., $\mathbf{L} \bar{\mathbf{W}} = \mathbf{V}$, where $\bar{\mathbf{W}} := \Sigma_{X_2 \bar{Y}}^\dagger \Sigma_{\bar{Y} Y}$. From our assumption that $\sigma_r(\Sigma_{\bar{Y} Y}^\dagger \Sigma_{\bar{Y} X_2}) = \beta$, we have $\|\bar{\mathbf{W}}\|_2 \leq \|\Sigma_{X_2 \bar{Y}}^\dagger \Sigma_{\bar{Y}}\|_2 \leq 1/\beta$. (Or we could directly define β as $\sigma_k(\Sigma_{\bar{Y} Y}^\dagger \Sigma_{\bar{Y} X_2}) \equiv \|\bar{\mathbf{W}}\|_2$.)

By concentration, we have $\mathbf{E} = \mathbf{X}_1 \Sigma_{X_1 X_1}^{-1} \Sigma_{X_1 X_2 | \bar{Y}}$ converges to $\Sigma_{X_1 X_1}^{-1/2} \Sigma_{X_1 X_2 | \bar{Y}}$. Specifically, when $n \gg k + \log 1/\delta$, $\|\mathbf{E}\|_F \leq 1.1 \|\Sigma_{X_1 X_1}^{-1/2} \Sigma_{X_1 X_2 | \bar{Y}}\|_F \leq 1.1 \epsilon_{\text{CI}}$ (by using Lemma C.2). Together we have $\|\mathbf{E} \bar{\mathbf{W}}\|_F \lesssim \epsilon_{\text{CI}}/\beta$.

Let $\hat{\mathbf{W}} = \arg \min_{\mathbf{W}} \|\mathbf{Y} - \Psi \mathbf{W}\|^2$. We note that $\mathbf{Y} = \mathbf{N} + \mathbf{V} = \mathbf{N} + \Psi \bar{\mathbf{W}} - \mathbf{E} \bar{\mathbf{W}}$ where \mathbf{V} is our target direction and \mathbf{N} is random noise (each row of \mathbf{N} has covariance matrix $\Sigma_{YY|X_1}$).

From basic inequality, we have:

$$\begin{aligned} \|\Psi \hat{\mathbf{W}} - \mathbf{Y}\|_F^2 &\leq \|\Psi \bar{\mathbf{W}} - \mathbf{Y}\|_F^2 = \|\mathbf{N} - \mathbf{E} \bar{\mathbf{W}}\|_F^2 \\ \Rightarrow \|\Psi \hat{\mathbf{W}} - \mathbf{V} - \mathbf{E} \bar{\mathbf{W}}\|^2 &\leq 2 \langle \Psi \hat{\mathbf{W}} - \mathbf{V} - \mathbf{E} \bar{\mathbf{W}}, \mathbf{N} - \mathbf{E} \bar{\mathbf{W}} \rangle \\ \Rightarrow \|\Psi \hat{\mathbf{W}} - \mathbf{V} - \mathbf{E} \bar{\mathbf{W}}\| &\leq \|P_{[\Psi, \mathbf{E}, \mathbf{V}]} \mathbf{N}\| + \|\mathbf{E} \bar{\mathbf{W}}\| \\ \Rightarrow \|\Psi \hat{\mathbf{W}} - \mathbf{V}\| &\lesssim \|\mathbf{E}\|_F \|\bar{\mathbf{W}}\| + (\sqrt{d_2} + \sqrt{\log 1/\delta}) \sqrt{\text{Tr}(\Sigma_{YY|X_1})} \\ &\quad \text{(from Lemma C.7)} \\ &\leq \sqrt{n_2} \frac{\epsilon_{\text{CI}}}{\beta} + (\sqrt{d_2} + \sqrt{\log 1/\delta}) \sqrt{\text{Tr}(\Sigma_{YY|X_1})} \\ &\quad \text{(from Assumption E.1)} \end{aligned}$$

Next, by the same procedure that concentrates $\frac{1}{n_2} \mathbf{X}_1^\top \mathbf{X}_1$ to $\Sigma_{X_1 X_1}$ with Claim C.2, we could easily get

$$\text{ER}[\hat{\mathbf{W}}] := \mathbb{E}[\|\hat{\mathbf{W}}^\top \psi(X_1) - f^*(X_1)\|^2] \lesssim \frac{\epsilon_{\text{CI}}^2}{\beta^2} + \text{Tr}(\Sigma_{YY|X_1}) \frac{d_2 + \log 1/\delta}{n_2}.$$

□

In the section below, we generalize the result from linear function space to arbitrary function space, and introduce the appropriate quantities to measure ACI.

E.2 Learnability with general function space

We state the main result with finite samples for both pretext task and downstream task to achieve good generalization. Let $\mathbf{X}_1^{\text{pre}} = [\mathbf{x}_1^{(1, \text{pre})}, \dots, \mathbf{x}_1^{(n_1, \text{pre})}]^\top \in \mathbb{R}^{n_1 \times d_1}$ and $\mathbf{X}_2 = [\mathbf{x}_2^{(1)}, \dots, \mathbf{x}_2^{(n_1)}]^\top \in \mathbb{R}^{n_1 \times d_2}$ be the training data from pretext task, where $(\mathbf{x}_1^{(i, \text{pre})}, \mathbf{x}_2^{(i)})$ is sampled from $P_{X_1 X_2}$. We consider two types of function spaces: $\mathcal{H} \in \{\mathcal{H}_1, \mathcal{H}_u\}$. Recall $\mathcal{H}_1 = \{\psi : \mathcal{X}_1 \rightarrow \mathbb{R}^{d_2} | \exists \mathbf{B} \in \mathbb{R}^{d_2 \times D_1}, \psi(\mathbf{x}_1) = \mathbf{B} \phi_1(\mathbf{x}_1)\}$ is induced by feature map $\phi_1 : \mathcal{X}_1 \rightarrow \mathbb{R}^{D_1}$. \mathcal{H}_u is a function space with universal approximation power (e.g. deep networks) that ensures $\psi^* = \mathbb{E}[X_2 | X_1] \in \mathcal{H}_u$. We learn a representation from \mathcal{H} by using n_1 samples: $\tilde{\psi} := \arg \min_{f \in \mathcal{H}_1^{d_2}} \frac{1}{n_1} \|\mathbf{X}_2 - f(\mathbf{X}_1^{\text{pre}})\|_F^2$. For

downstream tasks we similarly define $\mathbf{X}_1^{\text{down}} \in \mathbb{R}^{n_2 \times d_1}$, $\mathbf{Y} \in \mathbb{R}^{n_2 \times d_3^2}$, and learn a linear classifier trained on $\tilde{\psi}(\mathbf{X}_1^{\text{down}})$:

$$\hat{\mathbf{W}} \leftarrow \arg \min_{\mathbf{W}} \frac{1}{2n_2} \|\mathbf{Y} - \tilde{\psi}(\mathbf{X}_1^{\text{down}}) \mathbf{W}\|_F^2, \text{ER}_{\tilde{\psi}}(\hat{\mathbf{W}}) := \mathbb{E}_{X_1} \|f_{\mathcal{H}}^*(X_1) - \hat{\mathbf{W}}^\top \tilde{\psi}(X_1)\|_2^2.$$

Here $f_{\mathcal{H}}^* = \mathbb{E}^L[Y|\phi_1(X_1)]$ when $\mathcal{H} = \mathcal{H}_1$ and $f_{\mathcal{H}}^* = f^*$ for $\mathcal{H} = \mathcal{H}_u$.

Assumption E.2 (Correlation between X_2 and Y, Z). *Suppose there exists latent variable $Z \in \mathcal{Z}, |\mathcal{Z}| = m$ that ensures $\Sigma_{\phi_{\bar{y}} X_2}$ is full column rank and $\|\Sigma_{Y \phi_{\bar{y}}} \Sigma_{X_2 \phi_{\bar{y}}}^\dagger\|_2 = 1/\beta$, where A^\dagger is pseudo-inverse, and $\phi_{\bar{y}}$ is the one-hot embedding for $\bar{Y} = [Y, Z]$.*

Definition E.2 (Approximate conditional independence with function space \mathcal{H}).

1. For $\mathcal{H} = \mathcal{H}_1$, define $\epsilon_{CI} := \|\Sigma_{\phi_1 \phi_1}^{-1/2} \Sigma_{\phi_1 X_2 | \phi_{\bar{y}}}\|_F$.
2. For $\mathcal{H} = \mathcal{H}_u$, define $\epsilon_{CI}^2 := \mathbb{E}_{X_1} [\|\mathbb{E}[X_2|X_1] - \mathbb{E}_{\bar{Y}}[\mathbb{E}[X_2|\bar{Y}]|X_1]\|^2]$.

Exact CI for both cases ensures $\epsilon_{CI} = 0$. We present a unified analysis in the appendix that shows the ϵ_{CI} for the second case is same as the first case, with covariance operators instead of matrices.

When $\mathcal{H} = \mathcal{H}_u$, the residual term $N := Y - \mathbb{E}[Y|X_1]$ is mean zero and assumed to be σ^2 -subgaussian. When we use non-universal features ϕ_1 , $\mathbb{E}[Y - f_{\mathcal{H}_1}^*(X_1)|X_1]$ may not be mean zero. We thus introduce the standard assumption on $a := f^* - f_{\mathcal{H}_1}^* = \mathbb{E}[Y|X_1] - \mathbb{E}^L[Y|\phi_1(X_1)]$:

Assumption E.3. (Bounded approximation error [26]) *There exists a universal constant b_0 , such that $\|\Sigma_{\phi_1 \phi_1}^{-1/2} \phi_1(X_1) a(X_1)^\top\|_F \leq b_0 \sqrt{k}$ almost surely.*

Theorem E.3. *For a fixed $\delta \in (0, 1)$, under Assumptions E.2, E.3 for $\tilde{\psi}$ and ψ^* and 3.2 for non-universal feature maps, if $n_1, n_2 \gg \rho^4(d_2 + \log 1/\delta)$, and we learn the pretext tasks such that: $\mathbb{E} \|\tilde{\psi}(X_1) - \psi^*(X_1)\|_F^2 \leq \epsilon_{pre}^2$. Then the generalization error for downstream task with probability $1 - \delta$ is:*

$$\text{ER}_{\tilde{\psi}}(\hat{\mathbf{W}}) \leq \mathcal{O} \left(\sigma^2 \frac{d_2 + \log(d_2/\delta)}{n_2} + \frac{\epsilon_{CI}^2}{\beta^2} + \frac{\epsilon_{pre}^2}{\beta^2} \right) \quad (10)$$

We defer the proof to the appendix. The proof technique is similar to that of Section 3. The difference is now our $\tilde{\psi}(\mathbf{X}^{(\text{down})}) \in \mathbb{R}^{n_2 \times d_2}$ will be an approximately low rank matrix (low rank + small norm), where the low rank part is the high-signal features that implicitly comes from Y, Z that will be useful for downstream. The remaining part comes from ϵ_{CI} and ϵ_{pre} . Again by selecting the top km (dimension of $\phi_{\bar{y}}$) features we could further improve the sample complexity:

Remark E.1. *By applying PCA on $\tilde{\psi}(\mathbf{X}_1^{\text{down}})$ and keeping the top km principal components only, we can improve the bound in Theorem E.3 to*

$$\text{ER}_{\tilde{\psi}}(\hat{\mathbf{W}}) \leq \mathcal{O} \left(\sigma^2 \frac{km + \log(km/\delta)}{n_2} + \frac{\epsilon_{CI}^2}{\beta^2} + \frac{\epsilon_{pre}^2}{\beta^2} \right). \quad (11)$$

We take a closer look at the different sources of errors in (11): 1) the noise term $Y - f^*(X_1)$ with noise level σ^2 ; 2) ϵ_{CI} that measures the approximate CI; and 3) ϵ_{pre} the error from not learning the pretext task exactly. The first term is optimal setting ignoring log factors as we do linear regression on mk -dimensional features. The second and third term are non-reducible due to the fact that f^* is not exactly linear in ψ while we use it as a fixed feature and learn a linear function on it. Therefore it is important to fine-tune when we have sufficient downstream labeled data. We leave this as future work.

Compared to traditional supervised learning, learning $f_{\mathcal{H}}^*$ requires sample complexity scaling with the (Rademacher/Gaussian) complexity of \mathcal{H} (see e.g. [6, 44]), which is very large for complicated models such as deep networks.

In Section G, we consider a similar result for cross-entropy loss.

We leave the experiments to the appendix, where we verify our main Theorem (E.3) using simulations. We check that pretext task helps when CI is approximately satisfied in text domain, and demonstrate on a real-world image dataset that a pretext task-based linear model outperforms or is comparable to many baselines.

² $d_3 = k$ and $Y \equiv \phi_y(Y)$ (one-hot encoding) refers multi-class classification task, $d_3 = 1$ refers to regression.

F Omitted Proofs Beyond Conditional Independence

F.1 Technical Facts

Lemma F.1 (Approximation Error of PCA). *Let matrix $\mathbf{A} = \mathbf{L} + \mathbf{E}$ where \mathbf{L} is rank r < size of \mathbf{A} and $\|\mathbf{E}\|_2 \leq \epsilon$ and $\Sigma_r(\mathbf{A}) = \beta$. Then we have*

$$\|\sin \Theta(\mathbf{A}, \mathbf{L})\|_2 \leq \epsilon/\beta.$$

Proof. We use Davis Kahan for this proof. First note that $\|\mathbf{A} - \mathbf{L}\| = \|\mathbf{E}\| \leq \epsilon$. From Davis-Kahan we get:

$$\begin{aligned} \|\sin \Theta(\mathbf{A}, \mathbf{L})\|_2 &\leq \frac{\|\mathbf{E}\|_2}{\Sigma_r(\mathbf{A}) - \Sigma_{r+1}(\mathbf{L})} \\ &= \frac{\|\mathbf{E}\|_2}{\Sigma_r(\mathbf{A})} \\ &\lesssim \epsilon/\beta. \end{aligned}$$

□

F.2 Measuring conditional dependence with cross-covariance operator

In Definition E.2 we have two ways to quantify ACI based on the choices of \mathcal{H} . It is actually unified by the introduction of some cross-covariance operator norm. This subsection gives more details on it. $L^2(P_X)$ denotes the Hilbert space of square integrable function with respect to the measure P_X , the marginal distribution of X . We are interested in some function class $\mathcal{H}_x \subset L^2(P_X)$ that is induced from some feature maps:

Definition F.2 (General and Universal feature Map). *We denote feature map $\phi : \mathcal{X} \rightarrow \mathcal{F}$ that maps from a compact input space \mathcal{X} to the feature space \mathcal{F} . \mathcal{F} is a Hilbert space associated with inner product: $\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle_{\mathcal{F}}$. The associated function class is: $\mathcal{H}_x = \{h : \mathcal{X} \rightarrow \mathbb{R} | \exists w \in \mathcal{F}, h(\mathbf{x}) = \langle w, \phi(\mathbf{x}) \rangle_{\mathcal{F}}, \forall \mathbf{x} \in \mathcal{X}\}$. We call ϕ universal if the induced \mathcal{H}_x is dense in $L^2(P_X)$.*

Linear model is a special case when feature map $\phi = Id$ is identity mapping and the inner product is over Euclidean space. A feature map with higher order polynomials correspondingly incorporate high order moments [16, 20]. For discrete variable Y we overload ϕ as the one-hot embedding.

Remark F.1. *For continuous data, any universal kernel like Gaussian kernel or RBF kernel induce the universal feature map that we require [35]. Two-layer neural network with infinite width also satisfy it, i.e., $\forall \mathbf{x} \in \mathcal{X} \subset \mathbb{R}^d, \phi_{NN}(\mathbf{x}) : \mathcal{S}^{d-1} \times \mathbb{R} \rightarrow \mathbb{R}, \phi_{NN}(\mathbf{x})[\mathbf{w}, b] = \sigma(\mathbf{w}^\top \mathbf{x} + b)$ [5].*

When there's no ambiguity, we overload ϕ_1 as the random variable $\phi_1(X_1)$ over domain \mathcal{F}_1 , and \mathcal{H}_1 as the function class over X_1 . Next we characterize CI using the cross-covariance operator.

Definition F.3 (Cross-covariance operator). *For random variables $X \in \mathcal{X}, Y \in \mathcal{Y}$ with joint distribution $P : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$, and associated feature maps ϕ_x and ϕ_y , we denote by $\mathcal{C}_{\phi_x \phi_y} = \mathbb{E}[\phi_x(X) \otimes \phi_y(Y)] = \int_{\mathcal{X} \times \mathcal{Y}} \phi_x(x) \otimes \phi_y(y) dP(x, y)$, the (un-centered) cross-covariance operator. Similarly we denote by $\mathcal{C}_{X \phi_y} = \mathbb{E}[X \otimes \phi_y(Y)] : \mathcal{F}_y \rightarrow \mathcal{X}$.*

To understand what $\mathcal{C}_{\phi_x \phi_y}$ is, we note it is of the same shape as $\phi_x(x) \otimes \phi_y(y)$ for each individual $x \in \mathcal{X}, y \in \mathcal{Y}$. It can be viewed as a self-adjoint operator: $\mathcal{C}_{\phi_x \phi_y} : \mathcal{F}_y \rightarrow \mathcal{F}_x$, $\mathcal{C}_{\phi_x \phi_y} f = \int_{\mathcal{X} \times \mathcal{Y}} \langle \phi_y(y), f \rangle \phi_x(x) dP(x, y), \forall f \in \mathcal{F}_y$. For any $f \in \mathcal{H}_x$ and $g \in \mathcal{H}_y$, it satisfies: $\langle f, \mathcal{C}_{\phi_x \phi_y} g \rangle_{\mathcal{H}_x} = \mathbb{E}_{XY}[f(X)g(Y)]$ [4, 16]. CI ensures $\mathcal{C}_{\phi_1 X_2 | \phi_y} = 0$ for arbitrary ϕ_1, ϕ_2 :

Lemma F.4. *With one-hot encoding map ϕ_y and arbitrary $\phi_1, X_1 \perp X_2 | Y$ ensures:*

$$\mathcal{C}_{\phi_1 X_2 | \phi_y} := \mathcal{C}_{\phi_1 X_2} - \mathcal{C}_{\phi_1 \phi_y} \mathcal{C}_{\phi_y \phi_y}^{-1} \mathcal{C}_{\phi_y X_2} = 0. \quad (12)$$

A more complete discussion of cross-covariance operator and CI can be found in [16]. Also, recall that an operator $\mathcal{C} : \mathcal{F}_y \rightarrow \mathcal{F}_x$ is Hilbert-Schmidt (HS) [42] if for complete orthonormal systems (CONSs) $\{\zeta_i\}$ of \mathcal{F}_x and $\{\eta_i\}$ of \mathcal{F}_y , $\|\mathcal{C}\|_{\text{HS}}^2 := \sum_{i,j} \langle \zeta_j, \mathcal{C} \eta_i \rangle_{\mathcal{F}_x}^2 < \infty$. The Hilbert-Schmidt norm

627 generalizes the Frobenius norm from matrices to operators, and we will later use $\|\mathcal{C}_{\phi_1 X_2 | \phi_y}\|$ to
 628 quantify approximate CI.

629 We note that covariance operators [17, 16, 4] are commonly used to capture conditional dependence
 630 of random variables. In this work, we utilize the covariance operator to quantify the performance of
 631 the algorithm even when the algorithm is *not a kernel method*.

632 F.3 Omitted Proof in General Setting

633 **Claim F.5.** *For feature maps ϕ_1 with universal property, we have:*

$$\begin{aligned}\psi^*(X_1) &:= \mathbb{E}[X_2|X_1] = \mathbb{E}^L[X_2|\phi_1] \\ &= \mathcal{C}_{X_2 \phi_1} \mathcal{C}_{\phi_1 \phi_1}^{-1} \phi_1(X_1). \\ \text{Our target } f^*(X_1) &:= \mathbb{E}[Y|X_1] = \mathbb{E}^L[Y|\phi_1] \\ &= \mathcal{C}_{Y \phi_1} \mathcal{C}_{\phi_1 \phi_1}^{-1} \phi_1(X_1).\end{aligned}$$

634 *For general feature maps, we instead have:*

$$\begin{aligned}\psi^*(X_1) &:= \arg \min_{f \in \mathcal{H}_1^{d_2}} \mathbb{E}_{X_1 X_2} \|X_2 - f(X_1)\|_2^2 \\ &= \mathcal{C}_{X_2 \phi_1} \mathcal{C}_{\phi_1 \phi_1}^{-1} \phi_1(X_1). \\ \text{Our target } f^*(X_1) &:= \arg \min_{f \in \mathcal{H}_1^k} \mathbb{E}_{X_1 Y} \|Y - f(X_1)\|_2^2 \\ &= \mathcal{C}_{Y \phi_1} \mathcal{C}_{\phi_1 \phi_1}^{-1} \phi_1(X_1).\end{aligned}$$

635 To prove Claim F.5, we show the following lemma:

636 **Lemma F.6.** *Let $\phi : \mathcal{X} \rightarrow \mathcal{F}_x$ be a universal feature map, then for random variable $Y \in \mathcal{Y}$ we have:*

$$\mathbb{E}[Y|X] = \mathbb{E}^L[Y|\phi(X)].$$

637 *Proof of Lemma F.6.* Denote by $\mathbb{E}[Y|X = x] =: f(x)$. Since ϕ is dense in \mathcal{X} , there exists a linear
 638 operator $a : \mathcal{X} \rightarrow \mathbb{R}$ such that $\int_{x \in \mathcal{X}} a(x) \phi(x) [\cdot] dx = f(\cdot)$ a.e. Therefore the result comes directly
 639 from the universal property of ϕ . \square

640 *Proof of Claim F.5.* We want to show that for random variables Y, X , where X is associated with a
 641 universal feature map ϕ_x , we have $\mathbb{E}[Y|X] = \mathcal{C}_{Y \phi_x(X)} \mathcal{C}_{\phi_x(X) \phi_x(X)}^{-1} \phi_x(X)$.

642 First, from Lemma F.6, we have that $\mathbb{E}[Y|X] = \mathbb{E}^L[Y|\phi_x(X)]$. Next, write $A^* : \mathcal{F}_x \rightarrow \mathcal{Y}$ as the
 643 linear operator that satisfies

$$\begin{aligned}\mathbb{E}[Y|X] &= A^* \phi_x(X) \\ \text{s.t. } A^* &= \arg \min_A \mathbb{E}[\|Y - A \phi_x(X)\|^2].\end{aligned}$$

644 Therefore from the stationary condition we have $A^* \mathbb{E}_X[\phi_x(X) \otimes \phi_x(X)] = \mathbb{E}_{XY}[Y \otimes \phi_x(X)]$. Or
 645 namely we get $A^* = \mathcal{C}_{Y \phi_x} \mathcal{C}_{\phi_x \phi_x}^{-1}$ simply from the definition of the cross-covariance operator \mathcal{C} . \square

646 **Claim F.7.** $\|\mathcal{C}_{\phi_1 \phi_1}^{-1/2} \mathcal{C}_{\phi_1 X_2 | \phi_y}\|_{\text{HS}}^2 = \mathbb{E}_{X_1}[\|\mathbb{E}[X_2|X_1] - \mathbb{E}_{\bar{Y}}[\mathbb{E}[X_2|\bar{Y}]|X_1]\|^2] = \epsilon_{CI}^2$.

Proof.

$$\begin{aligned}&\|\mathcal{C}_{\phi_1 \phi_1}^{-1/2} \mathcal{C}_{\phi_1 X_2 | \phi_y}\|_{\text{HS}}^2 \\ &= \int_{X_1} \left\| \int_{X_2} \left(\frac{p_{X_1 X_2}(\mathbf{x}_1, \mathbf{x}_2)}{p_{X_1}(\mathbf{x}_1)} - \frac{p_{X_1 \perp X_2 | Y}(\mathbf{x}_1, \mathbf{x}_2)}{p_{X_1}(\mathbf{x}_1)} \right) X_2 dp_{\mathbf{x}_2} \right\|^2 dp_{\mathbf{x}_1} \\ &= \mathbb{E}_{X_1}[\|\mathbb{E}[X_2|X_1] - \mathbb{E}_{\bar{Y}}[\mathbb{E}[X_2|\bar{Y}]|X_1]\|^2].\end{aligned}$$

647 \square

648 F.4 Omitted Proof for Main Results

649 We first prove a simpler version without approximation error.

650 **Theorem F.8.** *For a fixed $\delta \in (0, 1)$, under Assumption E.2, 3.2, if there is no approximation error,*
 651 *i.e., there exists a linear operator A such that $f^*(X_1) \equiv A\phi_1(X_1)$, if $n_1, n_2 \gg \rho^4(d_2 + \log 1/\delta)$,*
 652 *and we learn the pretext tasks such that:*

$$\mathbb{E} \|\tilde{\psi}(X_1) - \psi^*(X_1)\|_F^2 \leq \epsilon_{pre}^2.$$

653 Then we are able to achieve generalization for downstream task with probability $1 - \delta$:

$$\mathbb{E}[\|f_{\mathcal{H}_1}^*(X_1) - \hat{\mathbf{W}}^\top \tilde{\psi}(X_1)\|^2] \leq \mathcal{O}\left\{\sigma^2 \frac{d_2 + \log d_2/\delta}{n_2} + \frac{\epsilon_{CI}^2}{\beta^2} + \frac{\epsilon_{pre}^2}{\beta^2}\right\}. \quad (13)$$

654 *Proof of Theorem F.8.* We follow the similar procedure as Theorem E.1. For the setting of no
 655 approximation error, we have $f^* = f_{\mathcal{H}_1}^*$, and the residual term $N := Y - f^*(X_1)$ is a mean-
 656 zero random variable with $\mathbb{E}[\|N\|^2|X_1] \lesssim \sigma^2$ according to our data assumption in Section 3.
 657 $\mathbf{N} = \mathbf{Y} - f^*(\mathbf{X}_1^{\text{down}})$ is the collected n_2 samples of noise terms. We write $Y \in \mathbb{R}^{d_3}$. For
 658 classification task, we have $Y \in \{e_i, i \in [k]\} \subset \mathbb{R}^k$ (i.e, $d_3 = k$) is one-hot encoded random variable.
 659 For regression problem, Y might be otherwise encoded. For instance, in the yearbook dataset, Y
 660 ranges from 1905 to 2013 and represents the years that the photos are taken. We want to note that our
 661 result is general for both cases: the bound doesn't depend on d_3 , but only depends on the variance of
 662 N .

663 Let $\Psi^*, \mathbf{L}, \mathbf{E}, \mathbf{V}$ be defined as follows:

664 Let $\mathbf{V} = f^*(\mathbf{X}_1^{\text{down}}) \equiv f_{\mathcal{H}_1}^*(\mathbf{X}_1^{\text{down}}) \equiv \phi(\mathbf{X}_1^{\text{down}})\mathcal{C}_{\phi_1}^{-1}\mathcal{C}_{\phi_1 Y}$ be our target direction. Denote the
 665 optimal representation matrix by

$$\begin{aligned} \Psi^* &:= \psi^*(\mathbf{X}_1^{\text{down}}) \\ &= \phi(\mathbf{X}_1^{\text{down}})\mathcal{C}_{\phi_1 \phi_1}^{-1}\mathcal{C}_{\phi_1 X_2} \\ &= \phi(\mathbf{X}_1^{\text{down}})\mathcal{C}_{\phi_1 \phi_1}^{-1}\mathcal{C}_{\phi_1 \phi_{\bar{y}}}\mathcal{C}_{\phi_{\bar{y}}}^{-1}\Sigma_{\phi_{\bar{y}} X_2} + \phi(\mathbf{X}_1^{\text{down}})\mathcal{C}_{\phi_1 \phi_1}^{-1}\mathcal{C}_{\phi_1 X_2|\phi_{\bar{y}}} \\ &=: \mathbf{L} + \mathbf{E}, \end{aligned}$$

666 where $\mathbf{L} = \phi(\mathbf{X}_1^{\text{down}})\mathcal{C}_{\phi_1 \phi_1}^{-1}\mathcal{C}_{\phi_1 \phi_{\bar{y}}}\mathcal{C}_{\phi_{\bar{y}}}^{-1}\mathcal{C}_{\phi_{\bar{y}} X_2}$ and $\mathbf{E} = \phi(\mathbf{X}_1^{\text{down}})\mathcal{C}_{\phi_1 \phi_1}^{-1}\mathcal{C}_{\phi_1 X_2|\bar{Y}}$.

667 In this proof, we denote S_Y as the matrix such that $S_Y \phi_{\bar{y}} = Y$. Specifically, if Y is of dimension d_3 ,
 668 S_Y is of size $d_3 \times |\mathcal{Y}||\mathcal{Z}|$. Therefore $S_Y \Sigma_{\phi_{\bar{y}} A} = \Sigma_{Y A}$ for any random variable A .

Therefore, similarly we have:

$$\mathbf{L}\Sigma_{X_2 \phi_{\bar{y}}}^\dagger \Sigma_{\phi_{\bar{y}} \phi_{\bar{y}}} S_Y^\top = \mathbf{L}\Sigma_{X_2 \phi_{\bar{y}}}^\dagger \Sigma_{\phi_{\bar{y}} Y} = \mathbf{L}\bar{\mathbf{W}} = \mathbf{V}$$

669 where $\bar{\mathbf{W}} := \Sigma_{X_2 \phi_{\bar{y}}}^\dagger \Sigma_{\phi_{\bar{y}} Y}$ satisfies $\|\bar{\mathbf{W}}\|_2 = 1/\beta$. Therefore $\text{span}(\mathbf{V}) \subseteq \text{span}(\mathbf{L})$ since we have
 670 assumed that $\Sigma_{X_2 \phi_{\bar{y}}}^\dagger \Sigma_{\phi_{\bar{y}} Y}$ to be full rank.

671 On the other hand, $\mathbf{E} = \mathbf{X}_1^{\text{down}}\mathcal{C}_{\phi_1 \phi_1}^{-1}\mathcal{C}_{\phi_1 X_2|\bar{Y}}$ concentrates to $\mathcal{C}_{\phi_1 \phi_1}^{-1/2}\mathcal{C}_{\phi_1 X_2|\phi_{\bar{y}}}$. Specifically, when
 672 $n \gg c + \log 1/\delta$, $\|\mathbf{E}\|_F \leq 1.1\|\mathcal{C}_{\phi_1 \phi_1}^{-1/2}\mathcal{C}_{\phi_1 X_2|\phi_{\bar{y}}}\|_F \leq 1.1\epsilon_{CI}$ (by using Lemma C.3). Together we
 673 have $\|\mathbf{E}\bar{\mathbf{W}}\|_F \lesssim \epsilon_{CI}/\beta$.

674 We also introduce the error from not learning ψ^* exactly: $\mathbf{E}^{\text{pre}} = \Psi - \Psi^* := \tilde{\psi}(\mathbf{X}_1^{\text{down}}) - \psi^*(\mathbf{X}_1^{\text{down}})$.
 675 With proper concentration and our assumption, we have that $\mathbb{E}\|\psi(X_1) - \psi^*(X_1)\|^2 \leq \epsilon_{pre}$ and
 676 $\frac{1}{\sqrt{n_2}}\|\psi(\mathbf{X}_1^{\text{down}}) - \psi^*(\mathbf{X}_1^{\text{down}})\|^2 \leq 1.1\epsilon_{pre}$.

677 Also, the noise term after projection satisfies $\|P_{[\Psi, \mathbf{E}, \mathbf{V}]} \mathbf{N}\| \lesssim \sqrt{d_2 + \log d_2/\delta}\sigma$ as using Lemma
 678 C.7. Therefore $\Psi = \Psi^* - \mathbf{E}^{\text{pre}} = \mathbf{L} + \mathbf{E} - \mathbf{E}^{\text{pre}}$.

679 Recall that $\hat{\mathbf{W}} = \arg \min_{\mathbf{W}} \|\psi(\mathbf{X}_1^{\text{down}})\mathbf{W} - \mathbf{Y}\|_F^2$. And with exactly the same procedure as Theorem
 680 E.1 we also get that:

$$\begin{aligned}\|\Psi\hat{\mathbf{W}} - \mathbf{V}\| &\leq 2\|\mathbf{E}\bar{\mathbf{W}}\| + 2\|\mathbf{E}^{\text{pre}}\bar{\mathbf{W}}\| + \|P_{[\Psi, \mathbf{E}, \mathbf{V}, \mathbf{E}^{\text{pre}}]}\mathbf{N}\| \\ &\lesssim \sqrt{n_2} \frac{\epsilon_{\text{CI}} + \epsilon_{\text{pre}}}{\beta} + \sigma \sqrt{d_2 + \log(d_2/\delta)}.\end{aligned}$$

681 With the proper concentration we also get:

$$\mathbb{E}[\|\hat{\mathbf{W}}^\top \psi(X_1) - f_{\mathcal{H}_1}^*(X_1)\|^2] \lesssim \frac{\epsilon_{\text{CI}}^2 + \epsilon_{\text{pre}}^2}{\beta^2} + \sigma^2 \frac{d_2 + \log(d_2/\delta)}{n_2}.$$

682

□

683 Next we move on to the proof of our main result Theorem E.3 where approximation error occurs.

684 *Proof of Theorem E.3.* The proof is a combination of Theorem 3.5 and Theorem F.8. We follow the
685 same notation as in Theorem F.8. Now the only difference is that an additional term $a(\mathbf{X}_1^{\text{down}})$ is
686 included in \mathbf{Y} :

$$\begin{aligned}\mathbf{Y} &= \mathbf{N} + f^*(\mathbf{X}_1^{\text{down}}) \\ &= \mathbf{N} + \Psi^* \bar{\mathbf{W}} + a(\mathbf{X}_1^{\text{down}}) \\ &= \mathbf{N} + (\Psi + \mathbf{E}^{\text{pre}}) \bar{\mathbf{W}} + a(\mathbf{X}_1^{\text{down}}) \\ &= \Psi \bar{\mathbf{W}} + (\mathbf{N} + \mathbf{E}^{\text{pre}} \bar{\mathbf{W}} + a(\mathbf{X}_1^{\text{down}})).\end{aligned}$$

687 From re-arranging $\frac{1}{2n_2} \|\mathbf{Y} - \Psi \hat{\mathbf{W}}\|_F^2 \leq \frac{1}{2n_2} \|\mathbf{Y} - \Psi \bar{\mathbf{W}}\|_F^2$,

$$\frac{1}{2n_2} \|\Psi(\bar{\mathbf{W}} - \hat{\mathbf{W}}) + (\mathbf{N} + \mathbf{E}^{\text{pre}} + a(\mathbf{X}_1^{\text{down}}))\|_F^2 \leq \frac{1}{2n_2} \|\mathbf{N} + \mathbf{E}^{\text{pre}} \bar{\mathbf{W}} + a(\mathbf{X}_1^{\text{down}})\|_F^2 \quad (14)$$

$$\Rightarrow \frac{1}{2n_2} \|\Psi(\bar{\mathbf{W}} - \hat{\mathbf{W}})\|_F^2 \leq \frac{1}{n_2} \langle \Psi(\bar{\mathbf{W}} - \hat{\mathbf{W}}), \mathbf{N} + \mathbf{E}^{\text{pre}} \bar{\mathbf{W}} + a(\mathbf{X}_1^{\text{down}}) \rangle. \quad (15)$$

688 Then with similar procedure as in the proof of Theorem 3.5, and write Ψ as $\phi(\mathbf{X}_1^{\text{down}})\mathbf{B}$, we have:

$$\begin{aligned}&\frac{1}{n_2} \langle \Psi(\bar{\mathbf{W}} - \hat{\mathbf{W}}), a(\mathbf{X}_1^{\text{down}}) \rangle \\ &= \frac{1}{n_2} \langle \mathbf{B}(\bar{\mathbf{W}} - \hat{\mathbf{W}}), \phi(\mathbf{X}_1^{\text{down}})^\top a(\mathbf{X}_1^{\text{down}}) \rangle \\ &= \frac{1}{n_2} \langle \mathcal{C}_{\phi_1}^{1/2} \mathbf{B}(\bar{\mathbf{W}} - \hat{\mathbf{W}}), \mathcal{C}_{\phi_1}^{-1/2} \phi(\mathbf{X}_1^{\text{down}})^\top a(\mathbf{X}_1^{\text{down}}) \rangle \\ &\leq \sqrt{\frac{d_2}{n_2}} \|\mathcal{C}_{\phi_1}^{1/2} \mathbf{B}(\bar{\mathbf{W}} - \hat{\mathbf{W}})\|_F \\ &\leq 1.1 \frac{1}{\sqrt{n_2}} \sqrt{\frac{d_2}{n_2}} \|\phi(\mathbf{X}_1^{\text{down}}) \mathbf{B}(\bar{\mathbf{W}} - \hat{\mathbf{W}})\|_F \\ &= 1.1 \frac{\sqrt{d_2}}{n_2} \|\Psi(\bar{\mathbf{W}} - \hat{\mathbf{W}})\|_F.\end{aligned}$$

689 Therefore plugging back to (15) we get:

$$\begin{aligned}&\frac{1}{2n_2} \|\Psi(\bar{\mathbf{W}} - \hat{\mathbf{W}})\|_F^2 \leq \frac{1}{n_2} \langle \Psi(\bar{\mathbf{W}} - \hat{\mathbf{W}}), \mathbf{N} + \mathbf{E}^{\text{pre}} \bar{\mathbf{W}} + a(\mathbf{X}_1^{\text{down}}) \rangle \\ &\Rightarrow \frac{1}{2n_2} \|\Psi(\bar{\mathbf{W}} - \hat{\mathbf{W}})\|_F \leq \frac{1}{2n_2} \|\mathbf{E}^{\text{pre}} \bar{\mathbf{W}}\|_F + \frac{1}{2n_2} \|P_\Psi \mathbf{N}\|_F + 1.1 \frac{\sqrt{d_2}}{n_2}. \\ &\Rightarrow \frac{1}{2\sqrt{n_2}} \|\Psi \hat{\mathbf{W}} - f_{\mathcal{H}_1}^*(\mathbf{X}_1^{\text{down}})\|_F - \|\mathbf{E} \bar{\mathbf{W}}\|_F \leq \frac{1}{\sqrt{n_2}} (1.1 \sqrt{d_2} + \|\mathbf{E}^{\text{pre}} \bar{\mathbf{W}}\| + \sqrt{d_2 + \log(d_2/\delta)}) \\ &\Rightarrow \frac{1}{2\sqrt{n_2}} \|\Psi \hat{\mathbf{W}} - f_{\mathcal{H}_1}^*(\mathbf{X}_1^{\text{down}})\|_F \lesssim \sqrt{\frac{d_2 + \log d_2/\delta}{n_2}} + \frac{\epsilon_{\text{CI}} + \epsilon_{\text{pre}}}{\beta}.\end{aligned}$$

690 Finally by concentrating $\frac{1}{n_2} \Psi^\top \Psi$ to $\mathbb{E}[\tilde{\psi}(X_1)\tilde{\psi}(X_1)^\top]$ we get:

$$\mathbb{E}[\|\hat{\mathbf{W}}^\top \tilde{\psi}(X_1) - f_{\mathcal{H}_1}^*(X_1)\|_2^2] \lesssim \frac{d_2 + \log d_2/\delta}{n_2} + \frac{\epsilon_{\text{CI}}^2 + \epsilon_{\text{pre}}^2}{\beta^2},$$

691 with probability $1 - \delta$. □

692 G Theoretical analysis for classification tasks

693 G.1 Classification tasks

694 We now consider the benefit of learning ψ from a class \mathcal{H}_1 on linear classification task for label set
695 $\mathcal{Y} = [k]$. The performance of a classifier is measured using the standard logistic loss

696 **Definition G.1.** For a task with $\mathcal{Y} = [k]$, classification loss for a predictor $f : \mathcal{X}_1 \rightarrow \mathbb{R}^k$ is

$$\ell_{\text{clf}}(f) = \mathbb{E}[\ell_{\log}(f(X_1), Y)], \text{ where } \ell_{\log}(\hat{y}, y) = \left[-\log \left(\frac{e^{\hat{y}_y}}{\sum_{y'} e^{\hat{y}_{y'}}} \right) \right]$$

697 The loss for representation $\psi : \mathcal{X}_1 \rightarrow \mathbb{R}^{d_1}$ and linear classifier $\mathbf{W} \in \mathbb{R}^{k \times d_1}$ is denoted by $\ell_{\text{clf}}(\mathbf{W}\psi)$.

698 We note that the function ℓ_{\log} is 1-Lipschitz in the first argument. The result will also hold for the
699 hinge loss $\ell_{\text{hinge}}(\hat{y}, y) = (1 - \hat{y}_y + \max_{y' \neq y} \hat{y}_{y'})_+$ which is also 1-Lipschitz, instead of ℓ_{\log} .

700 We assume that the optimal regressor $f_{\mathcal{H}_1}^*$ for one-hot encoding also does well on linear classification.
701

702 **Assumption G.1.** The best regressor for 1-hot encodings in \mathcal{H}_1 does well on classification, i.e.
703 $\ell_{\text{clf}}(\gamma f_{\mathcal{H}_1}^*) \leq \epsilon_{\text{one-hot}}$ is small for some scalar γ .

704 **Remark G.1.** Note that if \mathcal{H}_1 is universal, then $f_{\mathcal{H}_1}^*(\mathbf{x}_1) = \mathbb{E}[Y|X_1 = \mathbf{x}_1]$ and we know that $f_{\mathcal{H}_1}^*$
705 is the Bayes-optimal predictor for binary classification. In general one can potentially predict the
706 label by looking at $\arg \max_{i \in [k]} f_{\mathcal{H}_1}^*(\mathbf{x}_1)_i$. The scalar γ captures the margin in the predictor $f_{\mathcal{H}_1}^*$.

707 We now show that using the classifier $\hat{\mathbf{W}}$ obtained from linear regression on one-hot encoding with
708 learned representations $\tilde{\psi}$ will also be good on linear classification. The proof is in Section G

709 **Theorem G.2.** For a fixed $\delta \in (0, 1)$, under the same setting as Theorem E.3 and Assumption G.1,
710 we have:

$$\ell_{\text{clf}}(\gamma \hat{\mathbf{W}} \tilde{\psi}) \leq \mathcal{O} \left(\gamma \sqrt{\sigma^2 \frac{d_2 + \log d_2/\delta}{n_2} + \frac{\epsilon_{\text{CI}}^2}{\beta^2} + \frac{\epsilon_{\text{pre}}^2}{\beta^2}} \right) + \epsilon_{\text{one-hot}},$$

711 with probability $1 - \delta$.

712 **Proof of Theorem G.2.** We simply follow the following sequence of steps

$$\begin{aligned} \ell_{\text{clf}}(\gamma \hat{\mathbf{W}} \tilde{\psi}) &= \mathbb{E}[\ell_{\log}(\gamma \hat{\mathbf{W}} \tilde{\psi}(X_1), Y)] \\ &\stackrel{(a)}{\leq} \mathbb{E} \left[\ell_{\log}(\gamma f_{\mathcal{H}_1}^*(X_1), Y) + \gamma \|\hat{\mathbf{W}} \tilde{\psi}(X_1) - f_{\mathcal{H}_1}^*(X_1)\| \right] \\ &\stackrel{(b)}{\leq} \epsilon_{\text{one-hot}} + \gamma \sqrt{\mathbb{E} \left[\|\hat{\mathbf{W}} \tilde{\psi}(X_1) - f_{\mathcal{H}_1}^*(X_1)\|^2 \right]} \\ &= \epsilon_{\text{one-hot}} + \gamma \sqrt{\text{ER}_{\tilde{\psi}}[\hat{\mathbf{W}}]} \end{aligned}$$

713 where (a) follows because ℓ_{\log} is 1-Lipschitz and (b) follows from Assumption G.1 and Jensen's
714 inequality. Plugging in Theorem E.3 completes the proof. □

715 H Four Different Ways to Use CI

716 In this section we propose four different ways to use conditional independence to prove zero approxi-
717 mation error, i.e.,

718 **Claim H.1** (informal). *When conditional independence is satisfied: $X_1 \perp X_2 | Y$, and some non-*
719 *degeneracy is satisfied, there exists some matrix \mathbf{W} such that $\mathbb{E}[Y|X_1] = \mathbf{W} \mathbb{E}[X_2|X_1]$.*

720 We note that for simplicity, most of the results are presented for the jointly Gaussian case, where
721 everything could be captured by linear conditional expectation $\mathbb{E}^L[Y|X_1]$ or the covariance matrices. When generalizing the results for other random variables, we note just replace X_1, X_2, Y by
722 $\phi_1(X_1), \phi_2(X_2), \phi_y(Y)$ will suffice the same arguments.
723

724 H.1 Inverse Covariance Matrix

725 Write Σ as the covariance matrix for the joint distribution $P_{X_1 X_2 Y}$.

$$\Sigma = \begin{bmatrix} \Sigma_{XX} & \Sigma_{XY} \\ \Sigma_{YX}^\top & \Sigma_{YY} \end{bmatrix}, \quad \Sigma^{-1} = \begin{bmatrix} \mathbf{A} & \rho \\ \rho^\top & \mathbf{B} \end{bmatrix}$$

726 where $\mathbf{A} \in \mathbb{R}^{(d_1+d_2) \times (d_1+d_2)}, \rho \in \mathbb{R}^{(d_1+d_2) \times k}, \mathbf{B} \in \mathbb{R}^{k \times k}$. Furthermore

$$\rho = \begin{bmatrix} \rho_1 \\ \rho_2 \end{bmatrix}; \quad \mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix}$$

727 for $\rho_i \in \mathbb{R}^{d_i \times k}, i = 1, 2$ and $\mathbf{A}_{ij} \in \mathbb{R}^{d_i \times d_j}$ for $i, j \in \{1, 2\}$.

728 **Claim H.2.** *When conditional independence is satisfied, \mathbf{A} is block diagonal matrix, i.e., \mathbf{A}_{12} and*
729 *\mathbf{A}_{21} are zero matrices.*

730 **Lemma H.3.** *We have the following*

$$\mathbb{E}[X_1|X_2] = (\mathbf{A}_{11} - \bar{\rho}_1 \bar{\rho}_1^\top)^{-1} (\bar{\rho}_1 \bar{\rho}_2^\top - \mathbf{A}_{12}) X_2 \quad (16)$$

$$\mathbb{E}[X_2|X_1] = (\mathbf{A}_{22} - \bar{\rho}_2 \bar{\rho}_2^\top)^{-1} (\bar{\rho}_2 \bar{\rho}_1^\top - \mathbf{A}_{21}) X_1 \quad (17)$$

$$\mathbb{E}[Y|X] = -\mathbf{B}^{-\frac{1}{2}} (\bar{\rho}_1^\top X_1 + \bar{\rho}_2^\top X_2) \quad (18)$$

731 where $\bar{\rho}_i = \rho_i \mathbf{B}^{-\frac{1}{2}}$ for $i \in \{1, 2\}$. Also,

$$(\mathbf{A}_{11} - \bar{\rho}_1 \bar{\rho}_1^\top)^{-1} \bar{\rho}_1 \bar{\rho}_2^\top = \frac{1}{1 - \bar{\rho}_1^\top \mathbf{A}_{11}^{-1} \bar{\rho}_1} \mathbf{A}_{11}^{-1} \bar{\rho}_1 \bar{\rho}_2^\top$$

$$(\mathbf{A}_{22} - \bar{\rho}_2 \bar{\rho}_2^\top)^{-1} \bar{\rho}_2 \bar{\rho}_1^\top = \frac{1}{1 - \bar{\rho}_2^\top \mathbf{A}_{22}^{-1} \bar{\rho}_2} \mathbf{A}_{22}^{-1} \bar{\rho}_2 \bar{\rho}_1^\top$$

732 *Proof.* We know that $\mathbb{E}[X_1|X_2] = \Sigma_{12} \Sigma_{22}^{-1} X_2$ and $\mathbb{E}[X_2|X_1] = \Sigma_{21} \Sigma_{11}^{-1} x_1$, where

$$\Sigma_{XX} = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$$

733 First using $\Sigma \Sigma^{-1} = \mathbf{I}$, we get the following identities

$$\Sigma_{XX} \mathbf{A} + \Sigma_{XY} \rho^\top = \mathbf{I} \quad (19)$$

$$\Sigma_{XY}^\top \mathbf{A} + \Sigma_{YY} \rho^\top = 0 \quad (20)$$

$$\Sigma_{XX} \rho + \Sigma_{XY} \mathbf{B} = 0 \quad (21)$$

$$\Sigma_{XY}^\top \rho + \Sigma_{YY} \mathbf{B} = \mathbf{I} \quad (22)$$

734 From Equation (21) we get that $\Sigma_{XY} = -\Sigma_{XX} \rho \mathbf{B}^{-1}$ and plugging this into Equation (19) we get

$$\begin{aligned} & \Sigma_{XX} \mathbf{A} - \Sigma_{XX} \rho \mathbf{B}^{-1} \rho^\top = \mathbf{I} \\ \implies & \Sigma_{XX} = (\mathbf{A} - \rho \mathbf{B}^{-1} \rho^\top)^{-1} = (\mathbf{A} - \bar{\rho} \bar{\rho}^\top)^{-1} \\ \implies & \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} = \left(\begin{bmatrix} \mathbf{A}_{11} - \bar{\rho}_1 \bar{\rho}_1^\top & \mathbf{A}_{12} - \bar{\rho}_1 \bar{\rho}_2^\top \\ \mathbf{A}_{21} - \bar{\rho}_2 \bar{\rho}_1^\top & \mathbf{A}_{22} - \bar{\rho}_2 \bar{\rho}_2^\top \end{bmatrix} \right)^{-1} \end{aligned}$$

735 We now make use of the following expression for inverse of a matrix that uses Schur complement:
 736 $\mathbf{M}/\alpha = \delta - \gamma\alpha^{-1}\beta$ is the Schur complement of α for \mathbf{M} defined below

$$\text{If } \mathbf{M} = \begin{bmatrix} \alpha & \beta \\ \gamma & \delta \end{bmatrix}, \text{ then, } \mathbf{M}^{-1} = \begin{bmatrix} \alpha^{-1} + \alpha^{-1}\beta(\mathbf{M}/\alpha)^{-1}\gamma\alpha^{-1} & -\alpha^{-1}\beta(\mathbf{M}/\alpha)^{-1} \\ -(\mathbf{M}/\alpha)^{-1}\gamma\alpha^{-1} & (\mathbf{M}/\alpha)^{-1} \end{bmatrix}$$

737 For $\mathbf{M} = (\mathbf{A} - \bar{\rho}\bar{\rho}^\top)$, we have that $\Sigma_{XX} = \mathbf{M}^{-1}$ and thus

$$\begin{aligned} \Sigma_{12}\Sigma_{22}^{-1} &= -\alpha^{-1}\beta(\mathbf{M}/\alpha)^{-1}((\mathbf{M}/\alpha)^{-1})^{-1} \\ &= -\alpha^{-1}\beta \\ &= (\mathbf{A}_{11} - \bar{\rho}_1\bar{\rho}_1^\top)^{-1}(\bar{\rho}_1\bar{\rho}_2^\top - \mathbf{A}_{12}) \end{aligned}$$

738 This proves Equation (16) and similarly Equation (17) can be proved.

739 For Equation (18), we know that $\mathbb{E}[Y|X = (X_1, X_2)] = \Sigma_{YX}\Sigma_{XX}^{-1}X = \Sigma_{XY}^\top\Sigma_{XX}^{-1}X$. By using
 740 Equation (21) we get $\Sigma_{XY} = -\Sigma_{XX}\rho\mathbf{B}^{-1}$ and thus

$$\begin{aligned} \mathbb{E}[Y|X = (X_1, X_2)] &= -\mathbf{B}^{-1}\rho^\top\Sigma_{XX}\Sigma_{XX}^{-1}X \\ &= -\mathbf{B}^{-1}\rho^\top X = \mathbf{B}^{-1}(\rho_1^\top X_1 + \rho_2^\top X_2) \\ &= -\mathbf{B}^{-\frac{1}{2}}(\bar{\rho}_1^\top X_1 + \bar{\rho}_2^\top X_2) \end{aligned}$$

741 For the second part, we will use the fact that $(\mathbf{I} - \mathbf{a}\mathbf{b}^\top)^{-1} = \mathbf{I} + \frac{1}{1-\mathbf{a}^\top\mathbf{b}}\mathbf{a}\mathbf{b}^\top$. Thus

$$\begin{aligned} (\mathbf{A}_{11} - \bar{\rho}_1\bar{\rho}_1^\top)^{-1}\bar{\rho}_1\bar{\rho}_2 &= (\mathbf{I} - \mathbf{A}_{11}^{-1}\bar{\rho}_1\bar{\rho}_1^\top)\mathbf{A}_{11}^{-1}\bar{\rho}_1\bar{\rho}_2^\top \\ &= (\mathbf{I} + \frac{1}{1 - \bar{\rho}_1^\top\mathbf{A}_{11}^{-1}\bar{\rho}_1}\mathbf{A}_{11}^{-1}\bar{\rho}_1\bar{\rho}_1^\top)\mathbf{A}_{11}^{-1}\bar{\rho}_1\bar{\rho}_2^\top \\ &= \mathbf{A}_{11}^{-1}(\mathbf{I} + \frac{1}{1 - \bar{\rho}_1^\top\mathbf{A}_{11}^{-1}\bar{\rho}_1}\bar{\rho}_1\bar{\rho}_1^\top\mathbf{A}_{11}^{-1})\bar{\rho}_1\bar{\rho}_2^\top \\ &= \mathbf{A}_{11}^{-1}(\bar{\rho}_1\bar{\rho}_2^\top + \frac{\bar{\rho}_1\mathbf{A}_{11}^{-1}\bar{\rho}_1}{1 - \bar{\rho}_1^\top\mathbf{A}_{11}^{-1}\bar{\rho}_1}\bar{\rho}_1\bar{\rho}_2^\top) \\ &= \mathbf{A}_{11}^{-1}\bar{\rho}_1\bar{\rho}_2^\top(1 + \frac{\bar{\rho}_1\mathbf{A}_{11}^{-1}\bar{\rho}_1}{1 - \bar{\rho}_1^\top\mathbf{A}_{11}^{-1}\bar{\rho}_1}) \\ &= \frac{1}{1 - \bar{\rho}_1^\top\mathbf{A}_{11}^{-1}\bar{\rho}_1}\mathbf{A}_{11}^{-1}\bar{\rho}_1\bar{\rho}_2^\top \end{aligned}$$

742 The other statement can be proved similarly. □

Claim H.4.

$$\mathbb{E}[X_2|X_1] = (\mathbf{A}_{22} - \bar{\rho}_2\bar{\rho}_2^\top)^{-1}\bar{\rho}_2\bar{\rho}_1^\top X_1, \mathbb{E}[Y|X_1] = -\mathbf{B}^{-1/2}\bar{\rho}_1^\top X_1 - \mathbf{B}^{-1/2}\bar{\rho}_2^\top \mathbb{E}[X_2|X_1]$$

743 Therefore $\mathbb{E}[Y|X_1]$ is in the same direction as $\mathbb{E}[X_2|X_1]$.

744 H.2 Closed form of Linear Conditional Expectation

745 Refer to Claim B.1 and proof of Lemma B.2. As this is the simplest proof we used in our paper.

746 H.3 From Law of Iterated Expectation

$$\begin{aligned} \mathbb{E}^L[X_2|X_1] &= \mathbb{E}^L[\mathbb{E}^L[X_2|X_1, Y]|X_1] \\ &= \mathbb{E}\left[\left[\Sigma_{X_2X_1}, \Sigma_{X_2Y}\right]\left[\begin{matrix} \Sigma_{X_1X_1} & \Sigma_{X_1Y} \\ \Sigma_{YX_1} & \Sigma_{YY} \end{matrix}\right]^{-1}\begin{bmatrix} X_1 \\ Y \end{bmatrix} \middle| X_1\right] \\ &= \mathbf{A}X_1 + \mathbf{B}\mathbb{E}^L[Y|X_1]. \end{aligned}$$

747 Using block matrix inverse,

$$\begin{aligned} \mathbf{A} &= (\Sigma_{X_2 X_1} - \Sigma_{X_2 Y} \Sigma_{Y Y}^{-1} \Sigma_{Y X_1}) (\Sigma_{X_1 X_1} - \Sigma_{X_1 Y} \Sigma_{Y Y}^{-1} \Sigma_{Y X_1})^{-1} \in \mathbb{R}^{d_2 \times d_1} \\ &= \Sigma_{X_1 X_2 | Y} (\Sigma_{X_1 X_1 | Y})^{-1} \\ \mathbf{B} &= \Sigma_{X_2 Y | X_1} (\Sigma_{Y Y | X_1})^{-1} \in \mathbb{R}^{d_2 \times \mathcal{Y}}. \end{aligned}$$

748 Therefore in general (without conditional independence assumption) our learned representation will
749 be $\psi(x_1) = \mathbf{A}x_1 + \mathbf{B}f^*(x_1)$, where $f^*(\cdot) := \mathbb{E}^L[Y|X_1]$.

750 It's easy to see that to learn f^* from representation ψ , we need \mathbf{A} to have some good property, such
751 as light tail in eigenspace, and \mathbf{B} needs to be full rank in its column space.

752 Notice in the case of conditional independence, $\Sigma_{X_1 X_2 | Y} = 0$, and $\mathbf{A} = 0$. Therefore we could
753 easily learn f^* from ψ if X_2 has enough information of Y such that $\Sigma_{X_2 Y | X_1}$ is of the same rank as
754 dimension of Y .

755 **H.4 From $\mathbb{E}[X_2|X_1, Y] = \mathbb{E}[X_2|Y]$**

756 *Proof.* Let the representation function ψ be defined as follows, and let we use law of iterated
757 expectation:

$$\begin{aligned} \psi(\cdot) &:= \mathbb{E}[X_2|X_1] = \mathbb{E}[\mathbb{E}[X_2|X_1, Y]|X_1] \\ &= \mathbb{E}[\mathbb{E}[X_2|Y]|X_1] \quad (\text{uses CI}) \\ &= \sum_y P(Y = y|X_1) \mathbb{E}[X_2|Y = y] \\ &=: f(X_1)^\top \mathbf{A}, \end{aligned}$$

758 where $f : \mathbb{R}^{d_1} \rightarrow \Delta_{\mathcal{Y}}$ satisfies $f(x_1)_y = P(Y = y|X_1 = x_1)$, and $\mathbf{A} \in \mathbb{R}^{\mathcal{Y} \times d_2}$ satisfies $\mathbf{A}_{y,:} =$
759 $\mathbb{E}[X_2|Y = y]$. Here Δ_d denotes simplex of dimension d , which represents the discrete probability
760 density over support of size d .

761 Let $\mathbf{B} = \mathbf{A}^\dagger \in \mathbb{R}^{\mathcal{Y} \times d_2}$ be the pseudoinverse of matrix \mathbf{A} , and we get $\mathbf{B}\mathbf{A} = \mathbf{I}$ from our assumption
762 that \mathbf{A} is of rank $|\mathcal{Y}|$. Therefore $f(x_1) = \mathbf{B}\psi(x_1), \forall x_1$. Next we have:

$$\begin{aligned} \mathbb{E}[Y|X_1 = \mathbf{x}_1] &= \sum_y P(Y = y|X_1 = \mathbf{x}_1) \times y \\ &= \hat{\mathbf{Y}} f(\mathbf{x}_1) \\ &= (\hat{\mathbf{Y}} \mathbf{B}) \cdot \psi(X_1). \end{aligned}$$

763 Here we denote by $\hat{\mathbf{Y}} \in \mathbb{R}^{k \times \mathcal{Y}}$, $\hat{\mathbf{Y}}_{:,y} = y$ that spans the whole support \mathcal{Y} . Therefore let $\mathbf{W}^* = \hat{\mathbf{Y}} \mathbf{B}$
764 will finish the proof.

765 □

766 I Experiments

767 In this section, we empirically verify our claim that SSL performs well when ACI is satisfied.

768 **Simulations.** With synthetic data, we verify how excess risk (ER) scales with the cardinality/feature
769 dimension of \mathcal{Y} (k), and ACI (ϵ_{CI} in Definition E.2). We consider a mixture of Gaussian data and
770 conduct experiments with both linear function space (\mathcal{H}_1 with ϕ_1 as identity map) and universal
771 function space \mathcal{H}_u . We sample the label Y uniformly from $\{1, \dots, k\}$. For i -th class, the centers
772 $\mu_{1i} \in \mathbb{R}^{d_1}$ and $\mu_{2i} \in \mathbb{R}^{d_2}$ are uniformly sampled from $[0, 10]$. Given $Y = i$, $\alpha \in [0, 1]$, let
773 $X_1 \sim \mathcal{N}(\mu_{1i}, \mathbf{I})$, $\hat{X}_2 \sim \mathcal{N}(\mu_{2i}, \mathbf{I})$, and $X_2 = (1 - \alpha)\hat{X}_2 + \alpha X_1$. Therefore α is a correlation
774 coefficient: $\alpha = 0$ ensures X_2 being CI with X_1 given Y and when $\alpha = 1$, X_2 fully depends on X_1 .
775 (if $d_1 \neq d_2$, we append zeros or truncate to fit accordingly).

776 We first conduct experiments with linear function class. We learn a linear representation ψ with
777 n_1 samples and the linear prediction of Y from ψ with n_2 samples. We set $d_1 = 50$, $d_2 = 40$,

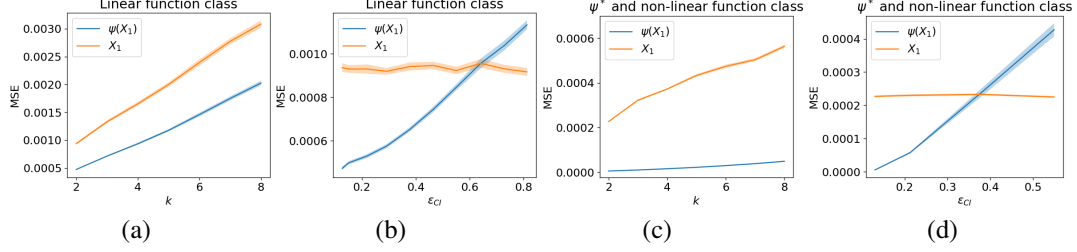


Figure 1: **Left two:** how MSE scales with k (the dimension of Y) and ϵ_{CI} (ACI E.2) with the linear function class. **Right two:** how MSE scales with k and ϵ with ψ^* and non-linear function class. Mean of 30 trials are shown in solid line and one standard error is shown by shadow.

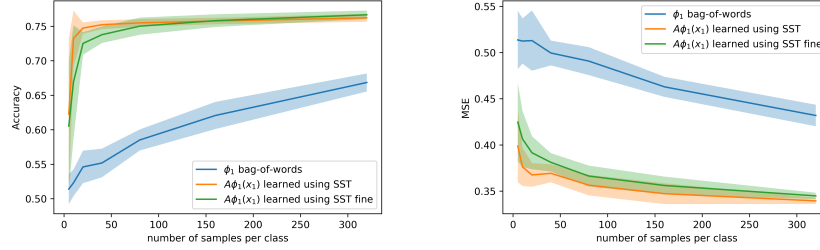


Figure 2: Performance on SST of baseline $\phi_1(x_1)$, i.e. bag-of-words, and learned $\psi(x_1)$ for the two settings. **Left:** Classification accuracy, **Right:** Regression MSE.

778 $n_1 = 4000$, $n_2 = 1000$ and ER is measured with Mean Squared Error (MSE). As shown in Figure
779 1(a)(b), the MSE of learning with $\psi(X_1)$ scales linearly with k as indicated in Theorem 3.5, and
780 scales linearly with ϵ_{CI} associated with linear function class as indicated in Theorem E.3. Next we
781 move on to general function class, i.e., $\psi^* = \mathbb{E}[Y|X_1]$ with a closed form solution (see example 3.1).
782 We use the same parameter settings as above. For baseline method, we use kernel linear regression to
783 predict Y using X_1 (we use RBF kernel which also has universal approximation power). As shown
784 in Figure 1(c)(d), the phenomenon is the same as what we observe in the linear function class setting,
785 and hence they respectively verify Theorem 3.2 and Theorem E.3 with \mathcal{H}_u .

786 **NLP task.** We look at the setting where both \mathcal{X}_1 and \mathcal{X}_2 are the set of sentences and perform
787 experiments by enforcing CI with and without latent variables. The downstream task is sentiment
788 analysis with the Stanford Sentiment Treebank (SST) dataset [45], where inputs are movie reviews
789 and the label set \mathcal{Y} is $\{\pm 1\}$. We use the representation class \mathcal{H}_1 , with features ϕ_1 being the bag-
790 of-words representation ($D_1 = 13848$). For X_2 we use a $d_2 = 300$ dimensional embedding of the
791 sentence, that is the mean of word vectors (random gaussians) for the words in the sentence. For
792 SSL data we consider 2 settings, (a) enforce CI with the labels \mathcal{Y} , (b) enforce CI with extra latent
793 variables, for which we use fine-grained version of SST with label set $\bar{\mathcal{Y}} = \{1, 2, 3, 4, 5\}^3$. We test
794 the learned ψ on SST binary task with linear regression and linear classification; results are presented
795 in Figure 2. We observe that in both settings ψ outperforms ϕ_1 , especially in the small-sample-size
796 regime. Also exact CI is better than CI with extra latent variables, as suggested by theory.

³Ratings $\{1, 2\}$ correspond to $y = -1$ and $\{4, 5\}$ correspond to $y = 1$

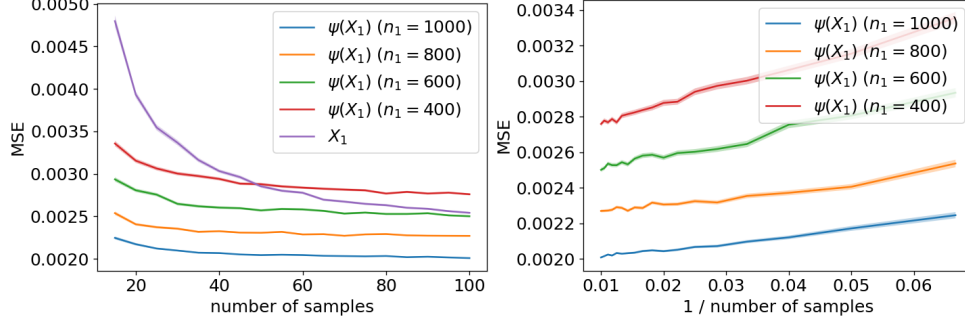


Figure 3: **Left:** MSE of using ψ to predict Y versus using X_1 directly to predict Y . Using ψ consistently outperforms using X_1 . **Right:** MSE of ψ learned with different n_1 . The MSE scale with $1/n_2$ as indicated by our analysis. Simulations are repeated 100 times, with the mean shown in solid line and one standard error shown in shadow.

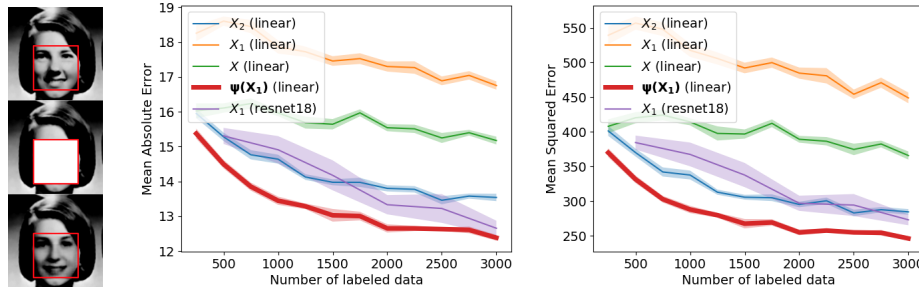


Figure 4: **Left:** Example of the X_2 (in the red box of the 1st row), the X_1 (out of the red box of the 1st row), the input to the inpainting task (the second row), $\psi(X_1)$ (the 3 row in the red box), and in this example $Y = 1967$. **Middle:** Mean Squared Error comparison of yearbook regression predicting dates. **Right:** Mean Absolute Error comparison of yearbook regression predicting dates. Experiments are repeated 10 times, with the mean shown in solid line and one standard error shown in shadow.

J More on the experiments

In this section, we describe more experiment results.

Simulations. Following Theorem E.3, we know that the Excessive Risk (ER) is also controlled by (1) the number of samples for the pretext task (n_1), and (2) the number of samples for the downstream task (n_2), besides k and ϵ_{CI} as discussed in the main text. In this simulation, we enforce strict conditional independence, and explore how ER varies with n_1 and n_2 . We generate the data the same way as in the main text, and keep $\alpha = 0$, $k = 2$, $d_1 = 50$ and $d_2 = 40$. We restrict the function class to linear model. Hence ψ is the linear model to predict X_2 from X_1 given the pretext dataset. We use Mean Squared Error (MSE) as the metric, since it is the empirical version of the ER. As shown in Figure 3, ψ consistently outperforms X_1 in predicting Y using a linear model learnt from the given downstream dataset, and ER does scale linearly with $1/n_2$, as indicated by our analysis.

Computer Vision Task. We testify if learning from ψ is more effective than learning directly from X_1 , in a realistic setting (without enforcing conditional independence). Specifically, we test on the Yearbook dataset [19], and try to predict the date when the portraits are taken (denoted as Y_D), which ranges from 1905 to 2013. We resize all the portraits to be 128 by 128. We crop out the center 64 by 64 pixels (the face), and treat it as X_2 , and treat the outer rim as X_1 as shown in Figure 4. Our task is to predict Y_D , which is the year when the portraits are taken, and the year ranges from 1905 to 2013. For ψ , we learn X_2 from X_1 with standard image inpainting techniques [40], and full set of training data (without labels). After that we fix the learned ψ and learn a linear model to predict Y_D from ψ

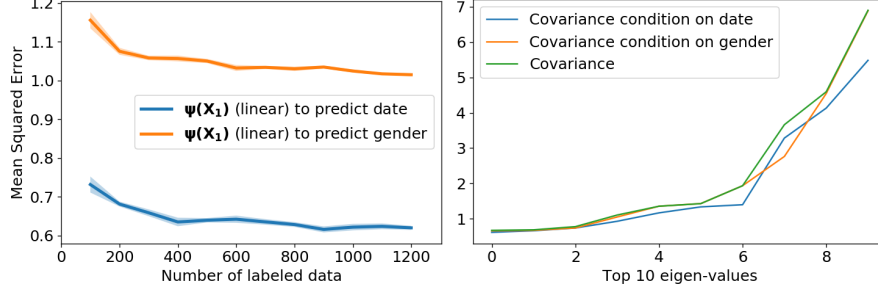


Figure 5: **Left:** Mean Squared Error comparison of predicting gender and predicting date. **Right:** the spectrum comparison of covariance condition on gender and condition on date.

816 using a smaller set of data (with labels). Besides linear model on X_1 , another strong baseline that
 817 we compare with is using ResNet18 [23] to predict Y_D from X_1 . With the full set of training data,
 818 this model is able to achieve a Mean Absolute Difference of 6.89, close to what state-of-the-art can
 819 achieve [19]. ResNet18 has similar amount of parameters as our generator, and hence roughly in the
 820 same function class. We show the MSE result as in Figure 4. Learning from ψ is more effective than
 821 learning from X_1 or X_2 directly, with linear model as well as with ResNet18. Practitioner usually
 822 fine-tune ψ with the downstream task, which usually leads to more competitive performance [40].

823 Following the same procedure, we try to predict the gender Y_G . We normalize the label (Y_G, Y_D) to
 824 unit variance, and confine ourself to linear function class. That is, instead of using a context encoder to
 825 impaint X_2 from X_1 , we confine ψ to be a linear function. As shown on the left of Figure 5, the MSE
 826 of predicting gender is higher than predicting dates. We find that $\|\Sigma_{\mathbf{X}_1\mathbf{X}_1}^{-1/2}\Sigma_{\mathbf{X}_1\mathbf{X}_2|Y_G}\|_F = 9.32$,
 827 while $\|\Sigma_{\mathbf{X}_1\mathbf{X}_1}^{-1/2}\Sigma_{\mathbf{X}_1\mathbf{X}_2|Y_D}\|_F = 8.15$. Moreover, as shown on the right of Figure 5, conditioning on
 828 Y_D cancels out more spectrum than conditioning on Y_G . In this case, we conjecture that, unlike Y_D ,
 829 Y_G does not capture much dependence between X_1 and X_2 . And as a result, ϵ_{CI} is larger, and the
 830 downstream performance is worse, as we expected.