

## A. Empirical Convergence of Vector Sequences

The LSL model in Section 2 and our main result in Section 4 require certain technical definitions.

**Definition 1** (Pseudo-Lipschitz continuity). For a given  $p \geq 1$ , a function  $\mathbf{f} : \mathbb{R}^d \rightarrow \mathbb{R}^m$  is called Pseudo-Lipschitz of order  $p$  if

$$\begin{aligned} \|\mathbf{f}(\mathbf{x}_1) - \mathbf{f}(\mathbf{x}_2)\| \\ \leq C \|\mathbf{x}_1 - \mathbf{x}_2\| (1 + \|\mathbf{x}_1\|^{p-1} + \|\mathbf{x}_2\|^{p-1}) \end{aligned} \quad (36)$$

for some constant  $C > 0$ .

Observe that for  $p = 1$ , the pseudo-Lipschitz is equivalent to the standard definition of Lipschitz continuity.

**Definition 2** (Uniform Lipschitz continuity). Let  $\phi(\mathbf{x}, \theta)$  be a function on  $\mathbf{r} \in \mathbb{R}^d$  and  $\theta \in \mathbb{R}^s$ . We say that  $\phi(\mathbf{x}, \theta)$  is *uniformly Lipschitz continuous* in  $\mathbf{x}$  at  $\theta = \bar{\theta}$  if there exists constants  $L_1, L_2 \geq 0$  and an open neighborhood  $U$  of  $\bar{\theta}$  such that

$$\|\phi(\mathbf{x}_1, \theta) - \phi(\mathbf{x}_2, \theta)\| \leq L_1 \|\mathbf{x}_1 - \mathbf{x}_2\| \quad (37)$$

for all  $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d$  and  $\theta \in U$ ; and

$$\|\phi(\mathbf{x}, \theta_1) - \phi(\mathbf{x}, \theta_2)\| \leq L_2 (1 + \|\mathbf{x}\|) \|\theta_1 - \theta_2\|, \quad (38)$$

for all  $\mathbf{x} \in \mathbb{R}^d$  and  $\theta_1, \theta_2 \in U$ .

**Definition 3** (Empirical convergence of a sequence). Consider a sequence of vectors  $\mathbf{x}(N) = \{\mathbf{x}_n(N)\}_{n=1}^N$  with  $\mathbf{x}_n(N) \in \mathbb{R}^d$ . So, each  $\mathbf{x}(N)$  is a block vector with a total of  $Nd$  components. For a finite  $p \geq 1$ , we say that the vector sequence  $\mathbf{x}(N)$  converges empirically with  $p$ -th order moments if there exists a random variable  $X \in \mathbb{R}^d$  such that

- (i)  $\mathbb{E}\|X\|_p^p < \infty$ ; and
- (ii) for any  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  that is pseudo-Lipschitz continuous of order  $p$ ,

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=1}^N f(\mathbf{x}_n(N)) = \mathbb{E}[f(X)]. \quad (39)$$

In this case, with some abuse of notation, we will write

$$\lim_{N \rightarrow \infty} \{\mathbf{x}_n\} \stackrel{PL(p)}{=} X, \quad (40)$$

where we have omitted the dependence on  $N$  in  $\mathbf{x}_n(N)$ . We note that the sequence  $\{\mathbf{x}(N)\}$  can be random or deterministic. If it is random, we will require that for every pseudo-Lipschitz function  $f(\cdot)$ , the limit (39) holds almost

surely. In particular, if  $\mathbf{x}_n \sim X$  are i.i.d. and  $\mathbb{E}\|X\|_p^p < \infty$ , then  $\mathbf{x}$  empirically converges to  $X$  with  $p^{\text{th}}$  order moments.

$PL(p)$  convergence is equivalent to weak convergence plus convergence in  $p$  moment (Bayati & Montanari, 2011), and hence  $PL(p)$  convergence is also equivalent to convergence in Wasserstein- $p$  metric (See Chapter 6. (Villani, 2008)). We use this fact later in proving Theorem 1.

## B. ML-VAMP Denoisers Details

Related to  $\mathbf{S}_{\text{mp}}$  and  $\mathbf{s}_{\text{mp}}$  from equation (11), we need to define two quantities  $\mathbf{s}_{\text{mp}}^+ \in \mathbb{R}^N$  and  $\mathbf{s}_{\text{mp}}^- \in \mathbb{R}^p$  that are zero-padded versions of the singular values  $\mathbf{s}_{\text{mp}}$ , so that for  $n > \min\{N, p\}$ , we set  $s_{\text{mp},n}^\pm = 0$ . Observe that  $(\mathbf{s}_{\text{mp}}^+)^2$  are eigenvalues of  $\mathbf{U}\mathbf{U}^T$  whereas  $(\mathbf{s}_{\text{mp}}^-)^2$  are eigenvalues of  $\mathbf{U}^T\mathbf{U}$ . Since  $\mathbf{s}_{\text{mp}}$  empirically converges to  $S_{\text{mp}}$  as given in (12), the vector  $\mathbf{s}_{\text{mp}}^+$  empirically converges to random variable  $S_{\text{mp}}^+$  whereas the vector  $\mathbf{s}_{\text{mp}}^-$  empirically converges to random variable  $S_{\text{mp}}^-$ , where a mass is placed at 0 appropriately. Specifically,  $S_{\text{mp}}^+$  has a point mass of  $(1 - \beta) + \delta_{\{0\}}$  when  $\beta < 1$ , whereas  $S_{\text{mp}}^-$  has a point mass of  $(1 - \frac{1}{\beta}) + \delta_{\{0\}}$ , when  $\beta > 1$ . In Appendix H (eqn. (113)), we provide the densities over positive parts of  $S_{\text{mp}}^+$  and  $S_{\text{mp}}^-$ .

A key property of our analysis will be that the non-linear functions (20) and the denoisers  $\mathbf{g}_\ell^\pm(\cdot)$  have simple forms.

Non-linear functions  $\phi_\ell(\cdot)$ : The non-linear functions all act *componentwise*. For example, for  $\phi_1(\cdot)$  in (20), we have

$$\mathbf{z}_1 = \phi_1(\mathbf{p}_0, \mathbf{s}_{\text{tr}}) = \text{diag}(\mathbf{s}_{\text{tr}})\mathbf{p}_0 \iff z_{1,n} = \phi_1(p_{0,n}, s_{\text{tr},n}),$$

where  $\phi_1(\cdot)$  is the scalar-valued function,

$$\phi_1(p_0, s) = sp_0. \quad (41)$$

Similarly, for  $\phi_2(\cdot)$ ,

$$\mathbf{z}_2 = \phi_2(\mathbf{p}_1, \mathbf{s}_{\text{mp}}^+) \iff z_{2,n} = \phi_2(\bar{p}_{1,n}, s_{\text{mp},n}^+), \quad n < N$$

where  $\bar{\mathbf{p}}_1 \in \mathbb{R}^N$  is the zero-padded version of  $\mathbf{p}_1$ , and

$$\phi_2(p_1, s) = sp_1. \quad (42)$$

Finally, the function  $\phi_3(\cdot)$  in (20) acts componentwise with

$$\phi_3(p_2, d) = \phi_{\text{out}}(p_2, d). \quad (43)$$

Input denoiser  $\mathbf{g}_0^+(\cdot)$ : Since  $F_0(\mathbf{z}_0) = F_{\text{in}}(\mathbf{z}_0)$ , and  $F_{\text{in}}(\cdot)$  given in (6), the denoiser (25a) acts *componentwise* in that,

$$\hat{\mathbf{z}}_0 = \mathbf{g}_0^+(\mathbf{r}_0^-, \gamma_0^-) \iff \hat{z}_{0,n} = g_0^+(r_{0,n}^-, \gamma_0^-),$$

where  $g_0^+(\cdot)$  is the scalar-valued function,

$$g_0^+(r_0^-, \gamma_0^-) := \underset{z_0}{\text{argmin}} f_{\text{in}}(z_0) + \frac{\gamma_0^-}{2} (z_0 - r_0^-)^2. \quad (44)$$

Thus, the vector optimization in (25a) reduces to a set of scalar optimizations (44) on each component.

**Output denoiser  $\mathbf{g}_3^-(\cdot)$ :** The output penalty  $F_3(\mathbf{p}_2, \mathbf{y}) = F_{\text{out}}(\mathbf{p}_2, \mathbf{y})$  where  $F_{\text{out}}(\mathbf{p}_2, \mathbf{y})$  has the separable form (6). Thus, similar to the case of  $\mathbf{g}_0(\cdot)$ , the denoiser  $\mathbf{g}_3(\cdot)$  in (25b) also acts componentwise with the function,

$$g_3^-(r_2^+, \gamma_2^+, y) := \operatorname{argmin}_{p_2} f_{\text{out}}(p_2, y) + \frac{\gamma_2^+}{2}(p_2 - r_2^+)^2. \quad (45)$$

**Linear denoiser  $\mathbf{g}_1^\pm(\cdot)$ :** The expressions for both denoisers  $g_1^\pm$  and  $g_2^\pm$  are very similar and can be explained together. The penalty  $F_1(\cdot)$  restricts  $\mathbf{z}_1 = \mathbf{S}_{\text{tr}}\mathbf{p}_0$ , where  $\mathbf{S}_{\text{tr}}$  is a square matrix. Hence, for  $\ell = 1$ , the minimization in (27) is given by,

$$\hat{\mathbf{p}}_0 := \operatorname{argmin}_{\mathbf{p}_0} \frac{\gamma_0^+}{2} \|\mathbf{p}_0 - \mathbf{r}_0^+\|^2 + \frac{\gamma_1^-}{2} \|\mathbf{S}_{\text{tr}}\mathbf{p}_0 - \mathbf{r}_1^-\|^2, \quad (46)$$

and  $\hat{\mathbf{z}}_1 = \mathbf{S}_{\text{tr}}\hat{\mathbf{p}}_0$ . This is a simple quadratic minimization and the components of  $\hat{\mathbf{p}}_0$  and  $\hat{\mathbf{z}}_1$  are given by

$$\begin{aligned} \hat{p}_{0,n} &= g_1^-(r_{0,n}^+, r_{1,n}^-, \gamma_0^+, \gamma_1^-, s_{\text{tr},n}) \\ \hat{z}_{1,n} &= g_1^+(r_{0,n}^+, r_{1,n}^-, \gamma_0^+, \gamma_1^-, s_{\text{tr},n}), \end{aligned}$$

where

$$g_1^-(r_0^+, r_1^-, \gamma_0^+, \gamma_1^-, s) := \frac{\gamma_0^+ r_0^+ + s\gamma_1^- r_1^-}{\gamma_0^+ + s^2\gamma_1^-} \quad (47a)$$

$$g_1^+(r_0^+, r_1^-, \gamma_0^+, \gamma_1^-, s) := \frac{s(\gamma_0^+ r_0^+ + s\gamma_1^- r_1^-)}{\gamma_0^+ + s^2\gamma_1^-} \quad (47b)$$

**Linear denoiser  $\mathbf{g}_2^\pm(\cdot)$ :** This denoiser is identical to the case  $\mathbf{g}_1^\pm(\cdot)$  in that we need to impose the linear constraint  $\mathbf{z}_2 = \mathbf{S}_{\text{mp}}\mathbf{p}_1$ . However  $\mathbf{S}_{\text{mp}}$  is in general a rectangular matrix and the two resulting cases of  $\beta \leq 1$  needs to be treated separately.

Recall the definitions of vectors  $\mathbf{s}_{\text{mp}}^+$  and  $\mathbf{s}_{\text{mp}}^-$  at the beginning of this section. Then, for  $\ell = 2$ , with the penalty  $F_2(\mathbf{p}_1, \mathbf{z}_2) = \delta_{\{\mathbf{z}_2 = \mathbf{S}_{\text{mp}}\mathbf{p}_1\}}$ , the solution to (27) has components,

$$\hat{p}_{1,n} = g_2^-(r_{1,n}^+, r_{2,n}^-, \gamma_1^+, \gamma_2^+, s_{\text{mp},n}^-) \quad (48a)$$

$$\hat{z}_{2,n} = g_2^+(r_{1,n}^+, r_{2,n}^-, \gamma_1^+, \gamma_2^+, s_{\text{mp},n}^+), \quad (48b)$$

with the identical functions  $g_2^- = g_1^-$  and  $g_2^+ = g_1^+$  as given by (47a) and (47b). Note that in (48a),  $n = 1, \dots, p$  and in (48b),  $n = 1, \dots, N$ .

## C. State Evolution Analysis of ML-VAMP

A key property of the ML-VAMP algorithm is that its performance in the LSL can be exactly described by a *scalar equivalent system*. In the scalar equivalent system, the vector-valued outputs of the algorithm are replaced by scalar random variables representing the typical behavior of the components of the vectors in the large-scale-limit (LSL). Each of the random variables are described by a set of parameters, where the parameters are given by a set of deterministic equations called the *state evolution* or SE.

The SE for the general ML-VAMP algorithm are derived in (Pandit et al., 2019) and the special case of the updates for ML-VAMP for GLM learning are shown in Algorithm 2 with details of functions  $\mathbf{g}_\ell^\pm$  in Appendix B. We see that the SE updates in Algorithm 2 parallel those in the ML-VAMP algorithm Algo. 1, except that vector quantities such as  $\hat{\mathbf{z}}_{k\ell}$ ,  $\hat{\mathbf{p}}_{k\ell}$ ,  $\mathbf{r}_{k\ell}^+$  and  $\mathbf{r}_{k\ell}^-$  are replaced by scalar random variables such as  $\hat{Z}_{k\ell}$ ,  $\hat{P}_{k\ell}$ ,  $R_{k\ell}^+$  and  $R_{k\ell}^-$ . Each of these random variables are described by the deterministic parameters such as  $\mathbf{K}_{k\ell} \in \mathbb{R}_{>0}^{2 \times 2}$ , and  $\tau_\ell^0, \tau_{k\ell}^- \in \mathbb{R}_+$ .

The updates in the section labeled as ‘‘Initial’’, provide the scalar equivalent model for the true system (18). In these updates,  $\Xi_\ell$  represent the limits of the vectors  $\xi_\ell$  in (19). That is,

$$\Xi_1 := S_{\text{tr}}, \quad \Xi_2 := S_{\text{mp}}^+, \quad \Xi_3 := D. \quad (49)$$

Due to assumptions in Section 2, we have that the components of  $\xi_\ell$  converge empirically as,

$$\lim_{N \rightarrow \infty} \{\xi_{\ell,i}\} \stackrel{PL(2)}{=} \Xi_\ell, \quad (50)$$

So, the  $\Xi_\ell$  represent the asymptotic distribution of the components of the vectors  $\xi_\ell$ .

The updates in sections labeled ‘‘Forward pass’’ and ‘‘Backward pass’’ in the SE equations in Algorithm 2 parallel those in Algorithm 1. The key quantities in these SE equations are the error variables,

$$\mathbf{p}_{k\ell}^+ := \mathbf{r}_{k\ell}^+ - \mathbf{p}_\ell^0, \quad \mathbf{q}_{k\ell}^- := \mathbf{r}_{k\ell}^- - \mathbf{z}_\ell^0,$$

which represent the errors of the estimates to the inputs of the denoisers. We will also be interested in their transforms,

$$\mathbf{q}_{k\ell}^+ = \mathbf{V}_\ell^\top \mathbf{p}_{k,\ell+1}^+, \quad \mathbf{p}_{k\ell}^- = \mathbf{V}_\ell \mathbf{q}_{k\ell}^-.$$

The following Theorem is an adapted version of the main result from (Pandit et al., 2019) to the iterates of Algorithms 1 and 2.

**Theorem 2.** *Consider the outputs of the ML-VAMP for GLM Learning Algorithm under the assumptions of Section 2. Assume the denoisers satisfy the continuity conditions in Assumption 1. Also, assume that the outputs of the SE satisfy*

$$\bar{\alpha}_{k\ell}^\pm \in (0, 1),$$

**Algorithm 2** SE for ML-VAMP for GLM Learning

```

1: // Initial
2: Initialize  $\bar{\gamma}_{0\ell} = \gamma_{0\ell}^-$  from Algorithm 1.
3:  $Q_{0\ell}^- \sim \mathcal{N}(0, \tau_{0\ell}^-)$  for some  $\tau_{0\ell}^- > 0$  for  $\ell = 0, 1, 2$ 
4:  $Z_0^0 = W^0$ 
5: for  $\ell = 0, \dots, L-1$  do
6:    $P_\ell^0 = \mathcal{N}(0, \tau_\ell^0)$ ,  $\tau_\ell^0 = \text{var}(Z_\ell^0)$ 
7:    $Z_{\ell+1}^0 = \phi_{\ell+1}(P_\ell^0, \Xi_{\ell+1})$ 
8: end for
9:
10: for  $k = 0, 1, \dots$  do
11:   // Forward Pass
12:   for  $\ell = 0, \dots, L-1$  do
13:     if  $\ell = 0$  then
14:        $R_{k0}^- = Z_\ell^0 + Q_{k0}^-$ 
15:        $\hat{Z}_{k0} = g_0^+(R_{k0}^-, \bar{\gamma}_{k0}^-)$ 
16:     else
17:        $R_{k,\ell-1}^+ = P_{\ell-1}^0 + P_{k,\ell-1}^+$ ,  $R_{k\ell}^- = Z_\ell^0 + Q_{k\ell}^-$ 
18:        $\hat{Z}_{k\ell} = g_\ell^+(R_{k,\ell-1}^+, R_{k\ell}^-, \bar{\gamma}_{k,\ell-1}^+, \bar{\gamma}_{k\ell}^-, \Xi_\ell)$ 
19:     end if
20:      $\bar{\alpha}_{k\ell}^+ = \mathbb{E} \partial \hat{Z}_{k\ell} / \partial Q_{k\ell}^-$ 
21:      $Q_{k\ell}^+ = \frac{\hat{Z}_{k\ell} - Z_\ell^0 - \bar{\alpha}_{k\ell}^+ Q_{k\ell}^-}{1 - \bar{\alpha}_{k\ell}^+}$ 
22:      $\bar{\gamma}_{k\ell}^+ = (\frac{1}{\bar{\alpha}_{k\ell}^+} - 1) \bar{\gamma}_{k\ell}^-$ 
23:      $(P_\ell^0, P_{k\ell}^+) \sim \mathcal{N}(0, \mathbf{K}_{k\ell}^+)$ ,  $\mathbf{K}_{k\ell}^+ = \text{cov}(Z_\ell^0, Q_{k\ell}^+)$ 
24:   end for
25:
26:   // Backward Pass
27:   for  $\ell = L, \dots, 1$  do
28:     if  $\ell = L$  then
29:        $R_{k,L-1}^+ = P_{L-1}^0 + P_{k,L-1}^+$ 
30:        $\hat{P}_{k,L-1} = g_L^-(R_{k,L-1}^+, \bar{\gamma}_{k,L-1}^+, Z_L^0)$ 
31:     else
32:        $R_{k,\ell-1}^+ = P_{\ell-1}^0 + P_{k,\ell-1}^+$ ,  $R_{k+1,\ell}^- = Z_\ell^0 + Q_{k+1,\ell}^-$ 
33:        $\hat{P}_{k,\ell-1} = g_\ell^-(R_{k,\ell-1}^+, R_{k+1,\ell}^-, \bar{\gamma}_{k,\ell-1}^+, \bar{\gamma}_{k+1,\ell}^-, \Xi_\ell)$ 
34:     end if
35:      $\bar{\alpha}_{k,\ell-1}^- = \mathbb{E} \partial \hat{P}_{k,\ell-1} / \partial P_{k,\ell-1}^+$ 
36:      $P_{k+1,\ell-1}^- = \frac{\hat{P}_{k,\ell-1} - P_{\ell-1}^0 - \bar{\alpha}_{k,\ell-1}^- P_{k,\ell-1}^+}{1 - \bar{\alpha}_{k,\ell-1}^-}$ 
37:      $\bar{\gamma}_{k+1,\ell-1}^- = (\frac{1}{\bar{\alpha}_{k,\ell-1}^-} - 1) \bar{\gamma}_{k,\ell-1}^+$ 
38:      $Q_{k+1,\ell-1}^- \sim \mathcal{N}(0, \tau_{k+1,\ell-1}^-)$ ,  $\tau_{k,\ell-1}^- = \mathbb{E}(P_{k+1,\ell-1}^-)^2$ 
39:   end for
40: end for

```

for all  $k$  and  $\ell$ . Suppose Algo. 1 is initialized so that the following convergence holds

$$\lim_{N \rightarrow \infty} \{\mathbf{r}_{0\ell}^- - \mathbf{z}_\ell^0\} \stackrel{PL(2)}{=} Q_{0\ell}^-$$

where  $(Q_{00}^-, Q_{01}^-, Q_{02}^-)$  are independent zero-mean Gaussians, independent of  $\{\Xi_\ell\}$ . Then,

(a) For any fixed iteration  $k \geq 0$  in the forward direction and layer  $\ell = 1, \dots, L-1$ , we have that, almost surely,

$$\lim_{N \rightarrow \infty} (\gamma_{k,\ell-1}^+, \gamma_{k\ell}^-) = (\bar{\gamma}_{k,\ell-1}^+, \bar{\gamma}_{k\ell}^-), \quad \text{and,} \quad (51)$$

$$\lim_{N \rightarrow \infty} \{(\hat{\mathbf{z}}_{k\ell}^+, \mathbf{z}_\ell^0, \mathbf{p}_{\ell-1}^0, \mathbf{r}_{k,\ell-1}^+, \mathbf{r}_\ell^-)\}$$

$$\stackrel{PL(2)}{=} (\hat{Z}_{k\ell}^+, Z_\ell^0, P_{\ell-1}^0, R_{k,\ell-1}^+, R_\ell^-) \quad (52)$$

where the variables on the right-hand side are from the SE equations (51) and (52) are the outputs of the SE equations in Algorithm 2. Similar equations hold for  $\ell = 0$  with the appropriate variables removed.

(b) Similarly, in the reverse direction, For any fixed iteration  $k \geq 0$  and layer  $\ell = 1, \dots, L-2$ , we have that, almost surely,

$$\lim_{N \rightarrow \infty} (\gamma_{k,\ell-1}^+, \gamma_{k+1,\ell}^-) = (\bar{\gamma}_{k,\ell-1}^+, \bar{\gamma}_{k+1,\ell}^-), \quad \text{and} \quad (53)$$

$$\lim_{N \rightarrow \infty} \{(\hat{\mathbf{p}}_{k+1,\ell-1}^+, \mathbf{z}_\ell^0, \mathbf{p}_{\ell-1}^0, \mathbf{r}_{k,\ell-1}^+, \mathbf{r}_{k+1,\ell}^-)\}$$

$$\stackrel{PL(2)}{=} (\hat{P}_{k+1,\ell-1}^+, Z_\ell^0, P_{\ell-1}^0, R_{k,\ell-1}^+, R_{k+1,\ell}^-). \quad (54)$$

Furthermore,  $(P_{\ell-1}^0, P_{k,\ell-1}^+)$  and  $Q_{k\ell}^-$  are independent.

*Proof.* This is a direct application of Theorem 3 from (Pandit et al., 2019) to the iterations of Algorithm 1. The convergence result in (Pandit et al., 2019) requires the uniform Lipschitz continuity of functions  $\mathbf{g}_\ell^\pm(\cdot)$ . Assumption 1 provides the required uniform Lipschitz continuity assumption on  $\mathbf{g}_0^+(\cdot)$  and  $\mathbf{g}_3^-(\cdot)$ . For the "middle" layers,  $\ell = 1, 2$ , the denoisers  $\mathbf{g}_\ell^\pm(\cdot)$  are linear and the uniform continuity assumption is valid since we have made the additional assumption that the terms  $\mathbf{s}_{\text{tr}}$  and  $\mathbf{s}_{\text{mp}}$  are bounded almost surely.  $\square$

A key use of the Theorem is to compute asymptotic empirical limits. Specifically, for a componentwise function  $\psi(\cdot)$ , let  $\langle \psi(\mathbf{x}) \rangle$  denotes the average  $\frac{1}{N} \sum_{n=1}^N \psi(x_n)$ . The above theorem then states that for any componentwise pseudo-Lipschitz function  $\psi(\cdot)$  of order 2, as  $N \rightarrow \infty$ , we have the following two properties

$$\lim_{N \rightarrow \infty} \langle \psi(\hat{\mathbf{z}}_{k\ell}^+, \mathbf{z}_\ell^0, \mathbf{p}_{\ell-1}^0, \mathbf{r}_{k,\ell-1}^+, \mathbf{r}_\ell^-) \rangle$$

$$= \mathbb{E} \psi(\hat{Z}_{k\ell}^+, Z_\ell^0, P_{\ell-1}^0, R_{k,\ell-1}^+, R_\ell^-)$$

$$\lim_{N \rightarrow \infty} \langle \psi(\hat{\mathbf{p}}_{k+1,\ell-1}^+, \mathbf{z}_\ell^0, \mathbf{p}_{\ell-1}^0, \mathbf{r}_{k,\ell-1}^+, \mathbf{r}_{k+1,\ell}^-) \rangle$$

$$= \mathbb{E} \psi(\hat{P}_{k+1,\ell-1}^+, Z_\ell^0, P_{\ell-1}^0, R_{k,\ell-1}^+, R_{k+1,\ell}^-).$$

That is, we can compute the empirical average over components with the expected value of the random variable limit. This convergence is key to proving Theorem 1.

## D. Empirical Convergence of Fixed Points

A consequence of Assumption 2 is that we can take the limit  $k \rightarrow \infty$  of the random variables in the SE algorithm. Specifically, let  $\mathbf{x}_k = \mathbf{x}_k(N)$  be any set of  $d$  outputs from the ML-VAMP for GLM Learning Algorithm under the assumptions of Theorem 2. Under Assumption 2, for each  $N$ , there exists a vector

$$\mathbf{x}(N) = \lim_{k \rightarrow \infty} \mathbf{x}_k(N), \quad (55)$$

representing the limit over  $k$ . For each  $k$ , Theorem 2 shows there also exists a random vector limit,

$$\lim_{N \rightarrow \infty} \{\mathbf{x}_{k,i}(N)\} \stackrel{PL(2)}{=} X_k, \quad (56)$$

representing the limit over  $N$ . The following proposition shows that we can take the limits of the random variables  $X_k$ .

**Proposition 1.** *Consider the outputs of the ML-VAMP for GLM Learning Algorithm under the assumptions of Theorem 2 and Assumption 2. Let  $\mathbf{x}_k = \mathbf{x}_k(N)$  be any set of  $d$  outputs from the algorithm and let  $\mathbf{x}(N)$  be its limit from (55) and let  $X_k$  be the random variable limit (56). Then, there exists a random variable  $X \in \mathbb{R}^d$  such that, for any pseudo-Lipschitz continuous  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ ,*

$$\lim_{k \rightarrow \infty} \mathbb{E}f(X_k) = \mathbb{E}f(X) = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N f(x_i(N)). \quad (57)$$

In addition, the SE parameter limits  $\bar{\alpha}_{k\ell}^\pm$  and  $\bar{\gamma}_{k\ell}^\pm$  converge to limits,

$$\lim_{k \rightarrow \infty} \bar{\alpha}_{k\ell}^\pm = \bar{\alpha}_\ell^\pm, \quad \lim_{k \rightarrow \infty} \bar{\gamma}_{k\ell}^\pm = \bar{\gamma}_\ell^\pm. \quad (58)$$

The proposition shows that, under the convergence assumption, Assumption 2, we can take the limits as  $k \rightarrow \infty$  of the random variables from the SE. To prove the proposition we first need the following simple lemma.

**Lemma 1.** *If  $\alpha_N$  and  $\beta_k \in \mathbb{R}$  are sequences such that*

$$\lim_{k \rightarrow \infty} \lim_{N \rightarrow \infty} |\alpha_N - \beta_k| = 0, \quad (59)$$

then, there exists a constant  $C$  such that,

$$\lim_{N \rightarrow \infty} \alpha_N = \lim_{k \rightarrow \infty} \beta_k = C. \quad (60)$$

In particular, the two limits in (60) exist.

*Proof.* For any  $\epsilon > 0$ , the limit (59) implies that there exists a  $k_\epsilon (\uparrow \infty$  as  $\epsilon \downarrow 0)$  such that for all  $k > k_\epsilon$ ,

$$\lim_{N \rightarrow \infty} |\alpha_N - \beta_k| < \epsilon,$$

from which we can conclude,

$$\liminf_{N \rightarrow \infty} \alpha_N > \beta_k - \epsilon$$

for all  $k > k_\epsilon$ . Therefore,

$$\liminf_{N \rightarrow \infty} \alpha_N \geq \sup_{k \geq k_\epsilon} \beta_k - \epsilon.$$

Since this is true for all  $\epsilon > 0$ , it follows that

$$\liminf_{N \rightarrow \infty} \alpha_N \geq \limsup_{k \rightarrow \infty} \beta_k. \quad (61)$$

Similarly,  $\limsup_{N \rightarrow \infty} \alpha_N \leq \inf_{k > k_\epsilon} \beta_k + \epsilon$ , whereby

$$\limsup_{N \rightarrow \infty} \alpha_N \leq \liminf_{k \rightarrow \infty} \beta_k. \quad (62)$$

Equations (61) and (62) together show that the limits in (60) exists and are equal.  $\square$

**Proof of Proposition 1** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be any pseudo-Lipschitz function of order 2, and define,

$$\alpha_N = \frac{1}{N} \sum_{i=1}^N f(x_i(N)), \quad \beta_k = \mathbb{E}f(X_k). \quad (63)$$

Their difference can be written as,

$$\alpha_N - \beta_k = A_{N,k} + B_{N,k}, \quad (64)$$

where

$$A_{N,k} := \frac{1}{N} \sum_{i=1}^N f(\mathbf{x}_i(N)) - f(\mathbf{x}_{k,i}(N)), \quad (65)$$

$$B_{N,k} := \frac{1}{N} \sum_{i=1}^N f(\mathbf{x}_{k,i}(N)) - \mathbb{E}f(X_k). \quad (66)$$

Since  $\{\mathbf{x}_{k,i}(N)\}$  converges PL(2) to  $X_k$ , we have,

$$\lim_{N \rightarrow \infty} B_{N,k} = 0. \quad (67)$$

For the term  $A_{N,k}$ ,

$$\begin{aligned} |A_{N,k}| &\stackrel{(a)}{\leq} \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N |f(\mathbf{x}_i(N)) - f(\mathbf{x}_{k,i}(N))| \\ &\stackrel{(b)}{\leq} \lim_{N \rightarrow \infty} \frac{C}{N} \sum_{i=1}^N a_{ki}(N)(1 + a_{ki}(N)) \\ &\stackrel{(c)}{\leq} C \lim_{N \rightarrow \infty} \sqrt{\frac{1}{N} \sum_{i=1}^N a_{ki}^2(N) + \frac{1}{N} \sum_{i=1}^N a_{ki}^2(N)} \\ &= C \lim_{N \rightarrow \infty} \epsilon_k(N)(1 + \epsilon_k(N)), \end{aligned} \quad (68)$$

where (a) follows from applying the triangle inequality to the definition of  $A_{N,k}$  in (65); (b) follows from the definition of pseudo-Lipschitz continuity in Definition 1,  $C > 0$  is the Lipschitz constant and

$$a_{ki}(N) := \|\mathbf{x}_{k,i}(N) - \mathbf{x}_i(N)\|_2,$$

and (c) follows from the RMS-AM inequality:

$$\left( \frac{1}{N} \sum_{i=1}^N a_{ki}(N) \right)^2 \leq \frac{1}{N} \sum_{i=1}^N a_{ki}^2(N) =: \epsilon_k^2(N).$$

By (29), we have that,

$$\lim_{k \rightarrow \infty} \lim_{N \rightarrow \infty} \epsilon_k(N) = 0.$$

Hence, from (68), it follows that,

$$\lim_{k \rightarrow \infty} \lim_{N \rightarrow \infty} A_{N,k} = 0. \quad (69)$$

Substituting (67) and (69) into (64) show that  $\alpha_N$  and  $\beta_k$  satisfy (59). Therefore, applying Lemma 1 we have that for any pseudo-Lipschitz function  $f(\cdot)$ , there exists a limit  $\Phi(f)$  such that,

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N f(x_i(N)) = \lim_{k \rightarrow \infty} \mathbb{E}f(X_k) = \Phi(f). \quad (70)$$

In particular, the two limits in (70) exists. When restricted to the continuous, bounded functions with the  $\|f\|_\infty$  norm, it is easy verified that  $\Phi(f)$  is a positive, linear, bounded function of  $f$ , with  $\Phi(1) = 1$ . Therefore, by the Riesz representation theorem, there exists a random variable  $X$  such that  $\Phi(f) = \mathbb{E}f(X)$ . This fact in combination with (70) shows (57).

It remains to prove the parameter limits in (58). We prove the result for the parameter  $\bar{\alpha}_{k\ell}^+$ . The proof for the other parameters are similar. Using Stein's lemma, it is shown in (Pandit et al., 2019) that

$$\bar{\alpha}_{k\ell}^+ = \frac{\mathbb{E}(\widehat{Z}_{k\ell} Q_{k\ell}^-)}{\mathbb{E}(Q_{k\ell}^-)^2}. \quad (71)$$

Since the numerator and denominator of (71) are  $PL(2)$  functions we have that the limit,

$$\begin{aligned} \bar{\alpha}_\ell^+ &:= \lim_{k \rightarrow \infty} \bar{\alpha}_{k\ell}^+ = \lim_{k \rightarrow \infty} \frac{\mathbb{E}(\widehat{Z}_{k\ell} Q_{k\ell}^-)}{\mathbb{E}(Q_{k\ell}^-)^2} \\ &= \frac{\mathbb{E}(\widehat{Z}_\ell Q_\ell^-)}{\mathbb{E}(Q_\ell^-)^2}, \end{aligned} \quad (72)$$

where  $\widehat{Z}_\ell$  and  $Q_\ell^-$  are the limits of  $\widehat{Z}_{k\ell}$  and  $Q_{k\ell}^-$ . This completes the proof.  $\square$

## E. Proof of Theorem 1

From Assumption 2, we know that for every  $N$ , every group of vectors  $\mathbf{x}_k$  converge to limits,  $\mathbf{x} := \lim_{k \rightarrow \infty} \mathbf{x}_k$ . The parameters,  $\gamma_{k\ell}^\pm$ , also converge to limits  $\bar{\gamma}_\ell^\pm := \lim_{k \rightarrow \infty} \gamma_{k\ell}^\pm$  for all  $\ell$ . By the continuity assumptions on the functions  $\mathbf{g}_\ell^\pm(\cdot)$ , the limits  $\mathbf{x}$  and  $\bar{\gamma}_\ell^\pm$  are fixed points of the algorithms.

A proof similar to that in (Pandit et al., 2019) shows that the fixed points  $\widehat{\mathbf{z}}_\ell$  and  $\widehat{\mathbf{p}}_\ell$  satisfy the KKT condition of the constrained optimization (22). This proves part (a).

The estimate  $\widehat{\mathbf{w}}$  is the limit,

$$\widehat{\mathbf{w}} = \widehat{\mathbf{z}}_0 = \lim_{k \rightarrow \infty} \widehat{\mathbf{z}}_{k0}.$$

Also, the true parameter is  $\mathbf{z}_0^0 = \mathbf{w}^0$ . By Proposition 1, we have that the  $PL(2)$  limits of these variables are

$$\lim_{N \rightarrow \infty} \{(\widehat{\mathbf{w}}, \mathbf{w}_0)\} \stackrel{PL(2)}{=} (\widehat{W}, W_0) := (\widehat{Z}_0, Z_0^0).$$

From line 2 of the SE Algorithm 2, we have

$$\widehat{W} = \widehat{Z}_0 = g_0^+(R_0^-, \bar{\gamma}_0^-) = \text{prox}_{f_{\text{in}}/\bar{\gamma}_0^-}(W^0 + Q_0^-).$$

This proves part (b).

To prove part (c), we use the limit

$$\lim_{N \rightarrow \infty} \{p_{0,n}^0, \widehat{p}_{0,n}\} \stackrel{PL(2)}{=} (P_0^0, \widehat{P}_0). \quad (73)$$

Since the fixed points are critical points of the constrained optimization (22),  $\widehat{\mathbf{p}}_0 = \mathbf{V}_0 \widehat{\mathbf{w}}$ . We also have  $\mathbf{p}_0^0 = \mathbf{V}_0 \mathbf{w}^0$ . Therefore,

$$\begin{aligned} [z_{\text{ts}}^{(N)} \widehat{z}_{\text{ts}}^{(N)}] &:= \mathbf{u}^\top \text{Diag}(\mathbf{s}_{\text{ts}}) \mathbf{V}_0 [\mathbf{w}^0 \widehat{\mathbf{w}}] \\ &= \mathbf{u}^\top \text{Diag}(\mathbf{s}_{\text{ts}}) [\mathbf{p}_0^0 \widehat{\mathbf{p}}_0]. \end{aligned} \quad (74)$$

Here,  $(N)$  in the subscript denotes the dependence on  $N$ . Since  $\mathbf{u} \sim \mathcal{N}(0, \frac{1}{p} \mathbf{I})$ ,  $[z_{\text{ts}}^{(N)} \widehat{z}_{\text{ts}}^{(N)}]$  is a zero-mean bivariate Gaussian with covariance matrix

$$\mathbf{M}^{(N)} = \frac{1}{p} \sum_{n=1}^p \begin{bmatrix} s_{\text{ts},n}^2 p_{0,n}^0 p_{0,n}^0 & s_{\text{ts},n}^2 p_{0,n}^0 \widehat{p}_{0,n} \\ s_{\text{ts},n}^2 p_{0,n}^0 \widehat{p}_{0,n} & s_{\text{ts},n}^2 \widehat{p}_{0,n} \widehat{p}_{0,n} \end{bmatrix}$$

The empirical convergence (73) yields the following limit,

$$\lim_{N \rightarrow \infty} \mathbf{M}^{(N)} = \mathbf{M} := \mathbb{E} S_{\text{ts}}^2 \begin{bmatrix} P_0^0 P_0^0 & P_0^0 \widehat{P}_0 \\ P_0^0 \widehat{P}_0 & \widehat{P}_0 \widehat{P}_0 \end{bmatrix}. \quad (75)$$

It suffices to show that the distribution of  $[z_{\text{ts}}^{(N)} \widehat{z}_{\text{ts}}^{(N)}]$  converges to the distribution of  $[Z_{\text{ts}} \widehat{Z}_{\text{ts}}]$  in the Wasserstein-2 metric as  $N \rightarrow \infty$ . (See the discussion in Appendix A on the equivalence of convergence in Wasserstein-2 metric and  $PL(2)$  convergence.)

Now, Wassestein-2 distance between two probability measures  $\nu_1$  and  $\nu_2$  is defined as

$$W_2(\nu_1, \nu_2) = \left( \inf_{\gamma \in \Gamma} \mathbb{E}_\gamma \|X_1 - X_2\|^2 \right)^{1/2}, \quad (76)$$

where  $\Gamma$  is the set of probability distributions on the product space with marginals consistent with  $\nu_1$  and  $\nu_2$ . For Gaussian measures  $\nu_1 = \mathcal{N}(\mathbf{0}, \Sigma_1)$  and  $\nu_2 = \mathcal{N}(\mathbf{0}, \Sigma_2)$  we have (Givens et al., 1984)

$$W_2^2(\nu_1, \nu_2) = \text{tr}(\Sigma_1 - 2(\Sigma_1^{1/2}\Sigma_2\Sigma_1^{1/2})^{1/2} + \Sigma_2).$$

Therefore, for Gaussian distributions  $\nu_1^{(N)} = \mathcal{N}(\mathbf{0}, \mathbf{M}^{(N)})$ , and  $\nu_2 = \mathcal{N}(\mathbf{0}, \mathbf{M})$ , the convergence (75) implies  $W_2(\nu_1^{(N)}, \nu_2) \rightarrow 0$ , i.e., convergence in Wasserstein-2 distance. Hence,

$$(z_{\text{ts}}^{(N)}, \hat{z}_{\text{ts}}^{(N)}) \xrightarrow{W_2} (Z_{\text{ts}}, \hat{Z}_{\text{ts}}) \sim \mathcal{N}(\mathbf{0}, \mathbf{M}),$$

where  $\mathbf{M}$  is the covariance matrix in (75). Hence the convergence holds in the PL(2) sense (see discussion in Appendix A on the equivalence of convergence in  $W_2$  and PL(2) convergence).

Hence the asymptotic generalization error (17) is

$$\begin{aligned} \mathcal{E}_{\text{ts}} &:= \lim_{N \rightarrow \infty} \mathbb{E} f_{\text{ts}}(\hat{y}_{\text{ts}}, y_{\text{ts}}) \\ &\stackrel{(a)}{=} \lim_{N \rightarrow \infty} \mathbb{E} f_{\text{ts}}(\phi_{\text{out}}(z_{\text{ts}}^{(N)}, D), \phi(\hat{z}_{\text{ts}}^{(N)})) \\ &\stackrel{(b)}{=} \mathbb{E} f_{\text{ts}}(\phi_{\text{out}}(Z_{\text{ts}}, D), \phi(\hat{Z}_{\text{ts}})), \end{aligned} \quad (77)$$

where (a) follows from (3); and step (b) follows from continuity assumption in Assumption 1(b) along with the definition of PL(2) convergence in Def. 3. This proves part (c).

## F. Formula for M

For the special cases in the next Appendix, it is useful to derive expressions for the entries the covariance matrix  $\mathbf{M}$  in (75). For the term  $m_{11}$ ,

$$m_{11} = \mathbb{E} S_{\text{ts}}^2 (P_0^0)^2 = \mathbb{E} S_{\text{ts}}^2 \mathbb{E} (P_0^0)^2 = \mathbb{E} S_{\text{ts}}^2 \cdot k_{11}, \quad (78)$$

where we have used the fact that  $P_0^0 \perp (S_{\text{ts}}, S_{\text{tr}})$ . Next,  $m_{12} = \mathbb{E} S_{\text{ts}}^2 P_0^0 \hat{P}_0$ . where,

$$\begin{aligned} \hat{P}_0 &= g_1^-(P_0^0 + P_0^+, Z_1^0 + Q_1^-, \bar{\gamma}_0^+, \bar{\gamma}_1^-, S_{\text{tr}}^-) \\ &= \frac{\bar{\gamma}_0^+ P_0^+ + S_{\text{tr}} \bar{\gamma}_1^- Q_1^-}{\bar{\gamma}_0^+ + S_{\text{tr}}^2 \bar{\gamma}_1^-} + P_0^0, \end{aligned} \quad (79)$$

where  $(P_0^0, P_0^+, Q_0^-)$  are independent of  $(S_{\text{tr}}, S_{\text{ts}})$ . Hence,

$$\begin{aligned} m_{12} &= \mathbb{E} S_{\text{ts}}^2 \cdot \mathbb{E} (P_0^0)^2 + \mathbb{E} \frac{S_{\text{ts}}^2 \bar{\gamma}_0^+}{\bar{\gamma}_0^+ + S_{\text{tr}}^2 \bar{\gamma}_1^-} \mathbb{E} [P_0^0 P_0^+] \\ &= m_{11} + \mathbb{E} \left( \frac{S_{\text{ts}}^2 \bar{\gamma}_0^+}{\bar{\gamma}_0^+ + S_{\text{tr}}^2 \bar{\gamma}_1^-} \right) \cdot k_{12}, \end{aligned} \quad (80)$$

since  $\mathbb{E}[P_0^0 Q_1^-] = 0$  and  $\mathbf{K}_0^+$  is the covariance matrix of  $(P_0^0, P_0^+)$  from line 2.

Finally for  $m_{22}$  we have,

$$\begin{aligned} m_{22} &= \mathbb{E} S_{\text{ts}}^2 \hat{P}_0 \hat{P}_0 \\ &= \mathbb{E} \left( \frac{S_{\text{ts}} \bar{\gamma}_0^+}{\bar{\gamma}_0^+ + S_{\text{tr}}^2 \bar{\gamma}_1^-} \right)^2 \mathbb{E} (P_0^+)^2 \\ &\quad + \mathbb{E} \left( \frac{S_{\text{ts}} S_{\text{tr}} \bar{\gamma}_1^-}{\bar{\gamma}_0^+ + S_{\text{tr}}^2 \bar{\gamma}_1^-} \right)^2 \mathbb{E} (Q_1^-)^2 \\ &\quad + \mathbb{E} S_{\text{ts}}^2 \mathbb{E} (P_0^0)^2 + 2 \mathbb{E} \frac{\bar{\gamma}_0^+ S_{\text{ts}}^2}{\bar{\gamma}_0^+ + \bar{\gamma}_1^- S_{\text{tr}}^2} \cdot \mathbb{E} P_0^0 P_0^+ \\ &= k_{22} \mathbb{E} \left( \frac{S_{\text{ts}} \bar{\gamma}_0^+}{\bar{\gamma}_0^+ + S_{\text{tr}}^2 \bar{\gamma}_1^-} \right)^2 \\ &\quad + \tau_1^- \mathbb{E} \left( \frac{S_{\text{ts}} S_{\text{tr}} \bar{\gamma}_1^-}{\bar{\gamma}_0^+ + S_{\text{tr}}^2 \bar{\gamma}_1^-} \right)^2 - m_{11} + 2m_{12}. \end{aligned} \quad (81)$$

## G. Special Cases

### G.1. Linear Output with Square Error

In this section we examine a few special cases of the GLM problem (2). We first consider a linear output with additive Gaussian noise and a squared error training and test loss. Specifically, consider the model,

$$\mathbf{y} = \mathbf{X} \mathbf{w}^0 + \mathbf{d} \quad (82)$$

We consider estimates of  $\mathbf{w}^0$  such that:

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\text{argmin}} \frac{1}{2} \|\mathbf{y} - \mathbf{X} \mathbf{w}\|^2 + \frac{\lambda}{2\beta} \|\mathbf{w}\|^2 \quad (83)$$

The factor  $\beta$  is added above since the two terms scale with a ratio of  $\beta$ . It does not change analysis. Consider the ML-VAMP GLM learning algorithm applied to this problem. The following corollary follows from the Main result in Theorem 1.

**Corollary 1** (Squared error). *For linear regression, i.e.,  $\phi(t) = t$ ,  $\phi_{\text{out}}(t, d) = t + d$ ,  $f_{\text{ts}}(y, \hat{y}) = (y_{\text{ts}} - \hat{y}_{\text{ts}})^2$ ,  $F_{\text{out}}(\mathbf{p}_2) = \frac{1}{N} \|\mathbf{y} - \mathbf{p}_2\|^2$ , we have*

$$\mathcal{E}_{\text{ts}}^{\text{LR}} = \mathbb{E} \left( \frac{\bar{\gamma}_0^+ S_{\text{ts}}}{\bar{\gamma}_0^+ + S_{\text{tr}}^2 \bar{\gamma}_1^-} \right)^2 k_{22} + \mathbb{E} \left( \frac{\bar{\gamma}_1^- S_{\text{tr}} S_{\text{ts}}}{\bar{\gamma}_0^+ + S_{\text{tr}}^2 \bar{\gamma}_1^-} \right)^2 \tau_1^- + \sigma_d^2.$$

*The quantities  $k_{22}$ ,  $\tau_1^-$ ,  $\bar{\gamma}_0^+$ ,  $\bar{\gamma}_1^-$  depend on the choice of regularizer  $\lambda$  and the covariance between features.*

*Proof.* This follows directly from the following observation:

$$\begin{aligned} \mathcal{E}_{\text{ts}}^{\text{SLR}} &= \mathbb{E} (Z_{\text{ts}} + D - \hat{Z}_{\text{ts}})^2 = \mathbb{E} (Z_{\text{ts}} - \hat{Z}_{\text{ts}})^2 + \mathbb{E} D^2 \\ &= m_{11} + m_{22} - 2m_{12} + \sigma_d^2. \end{aligned}$$

Substituting equation (81) proves the claim.  $\square$

## G.2. Ridge Regression with i.i.d. Covariates

We next the special case when the input features are independent, i.e., (83) where rows of  $\mathbf{X}$  corresponding to the training data has i.i.d Gaussian features with covariance  $\mathbf{P}_{\text{train}} = \frac{\sigma_{\text{tr}}^2}{p} \mathbf{I}$  and  $S_{\text{tr}} = \sigma_{\text{tr}}$ .

Although the solution to (83) exists in closed form  $(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$ , we can study the effect of the regularization parameter  $\lambda$  on the generalization error  $\mathcal{E}_{\text{ts}}$  as detailed in the result below.

**Corollary 2.** *Consider the ridge regression problem (83) with regularization parameter  $\lambda > 0$ . For the squared loss i.e.,  $f_{\text{ts}}(y, \hat{y}) = (y - \hat{y})^2$ , i.i.d Gaussian features without train-test mismatch, i.e.,  $S_{\text{tr}} = S_{\text{ts}} = \sigma_{\text{tr}}$ , the generalization error  $\mathcal{E}_{\text{ts}}^{\text{RR}}$  is given by Corollary 1, with constants*

$$k_{22} = \text{Var}(W^0), \quad \bar{\gamma}_0^+ = \lambda/\beta,$$

$$\bar{\gamma}_1^- = \begin{cases} \frac{1}{G} - \frac{\lambda}{\sigma_{\text{tr}}^2} & \beta < 1 \\ \frac{\frac{\lambda}{\sigma_{\text{tr}}^2 \beta} (\frac{1}{G} - \frac{\lambda}{\sigma_{\text{tr}}^2 \beta})}{\frac{\beta-1}{G} + \frac{\lambda}{\sigma_{\text{tr}}^2 \beta}} & \beta > 1 \end{cases}$$

where  $G = G_{\text{mp}}(-\frac{\lambda}{\sigma_{\text{tr}}^2 \beta})$ , with  $G_{\text{mp}}$  given in Appendix H, and  $\tau_1^- = \mathbb{E}(P_1^-)^2$  where  $P_1^-$  is given in equation (95) in the proof.

*Proof of Corollary 2.* We are interested in identifying the following constants appearing in Corollary 1:

$$\mathbf{K}_0^+, \tau_1^-, \bar{\gamma}_0^+, \bar{\gamma}_1^-.$$

These quantities are obtained as fixed points of the State Evolution Equations in Algo. 2. We explain below how to obtain expressions for these constants. Since these are fixed points we ignore the subscript  $k$  corresponding to the iteration number in Algo. 2.

In the case of problem (83), the maps  $\text{prox}_{f_{\text{in}}}$  and  $\text{prox}_{f_{\text{out}}}$ , i.e.,  $g_0^+$  and  $g_3^-$  respectively, can be expressed as closed-form formulae. This leads to simplification of the SE equations as explained below.

We start by looking at the *forward pass* (finding quantities with superscript '+' of Algorithm 2 for different layers and then the *backward pass* (finding quantities with superscript '-' to get the parameters  $\{\mathbf{K}_\ell^+, \tau_\ell^-, \bar{\alpha}_\ell^\pm, \bar{\gamma}_\ell^\pm\}$  for  $\ell = 0, 1, 2$ ).

To begin with, notice that  $f_{\text{in}}(w) = \frac{\lambda}{2} w^2$ , and therefore the denoiser  $g_0^+(\cdot)$  in (44) is simply,

$$g_0^+(r_0^-, \gamma_0^-) = \frac{\gamma_0^-}{\gamma_0^- + \lambda/\beta} r_0^-, \quad \text{and} \quad \frac{\partial g_0^+}{\partial r_0^-} = \frac{\gamma_0^-}{\gamma_0^- + \lambda/\beta}$$

Using the random variable  $R_0^-$  and substituting in the expression of the denoiser to get  $\hat{Z}_0$ , we can now calculate  $\bar{\alpha}_0^+$

using lines 2 and 2,

$$\bar{\alpha}_0^+ = \frac{\bar{\gamma}_0^-}{\bar{\gamma}_0^- + \lambda/\beta}, \quad \bar{\gamma}_0^+ = \lambda/\beta. \quad (84)$$

Similarly, we have  $f_{\text{out}}(p_2) = \frac{1}{2}(p_2 - y)^2$ , whereby the output denoiser  $g_3^-(\cdot)$  in the last layer for ridge regression is given by,

$$g_3^-(r_2^+, \gamma_2^+, y) = \frac{\gamma_2^+ r_2^+ + y}{\gamma_2^+ + 1}. \quad (85)$$

By substituting this denoiser in line 2 of the algorithm we get  $\hat{P}_2^-$  and thus, following the lines 2-2 of the algorithm we have

$$\bar{\alpha}_2^- = \frac{\bar{\gamma}_2^+}{\bar{\gamma}_2^+ + 1}, \quad \text{whereby} \quad \bar{\gamma}_2^- = 1. \quad (86)$$

Having identified these constants  $\bar{\alpha}_0^+, \bar{\gamma}_0^+, \bar{\alpha}_2^-, \bar{\gamma}_2^-$ , we will now sequentially identify the quantities

$$(\bar{\alpha}_0^+, \bar{\gamma}_0^+) \rightarrow \mathbf{K}_0^+ \rightarrow (\bar{\alpha}_1^-, \bar{\gamma}_1^-) \rightarrow \mathbf{K}_1^+ \rightarrow (\bar{\alpha}_2^-, \bar{\gamma}_2^-) \rightarrow \mathbf{K}_2^+$$

in the forward pass, and then the quantities

$$\tau_0^- \leftarrow (\bar{\alpha}_0^-, \bar{\gamma}_0^-) \leftarrow \tau_1^- \leftarrow (\bar{\alpha}_1^-, \bar{\gamma}_1^-) \leftarrow \tau_2^- \leftarrow (\bar{\alpha}_2^-, \bar{\gamma}_2^-)$$

in the backward pass.

We also note that we have

$$\bar{\alpha}_\ell^+ + \bar{\alpha}_\ell^- = 1 \quad (87)$$

**Forward Pass:** Observe that  $\mathbf{K}_0^+ = \text{Cov}(Z_0, Q_0^+)$ . Now, from line 2, on simplification we get  $Q_0^+ = -W_0^0$  whereby,

$$\mathbf{K}_0^+ = \text{var}(W^0) \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}. \quad (88)$$

Notice that from line 2, the pair  $(P_0^0, P_0^+)$  is jointly Gaussian with covariance matrix  $\mathbf{K}_0^+$ . But the above equation means that  $P_0^+ = -P_0^0$ , whereby  $R_0^+ = 0$  from line 2.

Now, the linear denoiser  $g_1^+(\cdot)$  is defined as in (47a). Note that since we are considering i.i.d Gaussian features for this problem, the random variable  $S_{\text{tr}}$  in this layer is a constant  $\sigma_{\text{tr}}$ . Therefore, similar to layer  $\ell = 0$  by evaluating lines 2-2 of the algorithm we get  $Q_1^+ = -Z_1^0$ , whereby

$$\bar{\alpha}_1^+ = \frac{\sigma_{\text{tr}}^2 \bar{\gamma}_1^-}{\bar{\gamma}_0^+ + \sigma_{\text{tr}}^2 \bar{\gamma}_1^-}, \quad \bar{\gamma}_1^+ = \frac{\bar{\gamma}_0^+}{\sigma_{\text{tr}}^2} = \frac{\lambda}{\sigma_{\text{tr}}^2 \beta}, \quad \mathbf{K}_1^+ = \sigma_{\text{tr}}^2 \mathbf{K}_0^+. \quad (89)$$

Observe that this means

$$P_1^+ = -P_1^0. \quad (90)$$

**Backward Pass:** Since  $Y = \phi_{\text{out}}(P_2^0, D) = P_2^0 + D$ , line 2 of algorithm on simplification yields  $P_2^- = D$ , whereby we can get  $\tau_2^-$ ,

$$\tau_2^- = \mathbb{E}(P_2^-)^2 = \mathbb{E}[D^2] = \sigma_d^2. \quad (91)$$

Next, to calculate the terms  $(\bar{\alpha}_1^-, \bar{\gamma}_1^-)$ , we use the denoiser  $g_2^-$  defined in (47a) for line 2 of the algorithm to get  $\hat{P}_1$ .

$$\hat{P}_1 = \frac{\bar{\gamma}_1^+ R_1^+ + S_{\text{mp}}^- \bar{\gamma}_2^- R_2^-}{\bar{\gamma}_1^+ + (S_{\text{mp}}^-)^2 \bar{\gamma}_2^-} = \frac{S_{\text{mp}}^- (S_{\text{mp}}^+ P_1^0 + Q_2^-)}{\bar{\gamma}_1^+ + (S_{\text{mp}}^-)^2}, \quad (92)$$

where we have used  $\bar{\gamma}_2^- = 1$ ,  $R_1^+ = P_1^0 + P_1^+ = 0$  due to (90), and  $R_2^- = Z_2^0 + Q_2^- = S_{\text{mp}}^+ P_1^0 + Q_2^-$  from lines 2, 2 and 2 respectively.

Then, we calculate  $\bar{\alpha}_1^-$  and  $\bar{\gamma}_1^-$  as  $\bar{\alpha}_1^- = \mathbb{E} \frac{\partial g_2^-}{\partial P_1^+} = \mathbb{E} \frac{\bar{\gamma}_1^+}{\bar{\gamma}_1^+ + (S_{\text{mp}}^-)^2}$ . This gives,

$$\bar{\alpha}_1^- = \begin{cases} \frac{\lambda}{\sigma_{\text{tr}}^2 \beta} G & \beta < 1 \\ (1 - \frac{1}{\beta}) + \frac{1}{\beta} \frac{\lambda}{\sigma_{\text{tr}}^2 \beta} G & \beta \geq 1 \end{cases} \quad (93)$$

Here, in the overparameterized case ( $\beta > 1$ ), the denoiser  $g_2^-$  outputs  $R_1^+$  with probability  $1 - \frac{1}{\beta}$  and  $\frac{\lambda}{\sigma_{\text{tr}}^2 \beta} G$  with probability  $\frac{1}{\beta}$ .

Next, from line 2 we get,

$$\bar{\gamma}_1^- = (\frac{1}{\bar{\alpha}_1^-} - 1) \bar{\gamma}_1^+ = \begin{cases} \frac{1}{G} - \frac{\lambda}{\sigma_{\text{tr}}^2 \beta} & \beta < 1 \\ \frac{\sigma_{\text{tr}}^2 \beta (\frac{1}{G} - \frac{\lambda}{\sigma_{\text{tr}}^2 \beta})}{\frac{\beta-1}{G} + \frac{\lambda}{\sigma_{\text{tr}}^2 \beta}} & \beta > 1 \end{cases} \quad (94)$$

Now from line 2 and equation (87) we get,

$$\begin{aligned} \bar{\alpha}_1^+ P_1^- &= \hat{P}_1 - P_1^0 - \bar{\alpha}_1^- P_1^+ \stackrel{(a)}{=} \hat{P}_1 - \bar{\alpha}_1^+ P_1^0 \\ &\stackrel{(b)}{=} \underbrace{\left( \frac{S_{\text{mp}}^- S_{\text{mp}}^+}{\sigma_{\text{tr}}^2 \beta + (S_{\text{mp}}^-)^2} - \bar{\alpha}_1^+ \right)}_A P_1^0 + \underbrace{\frac{S_{\text{mp}}^-}{\sigma_{\text{tr}}^2 \beta + (S_{\text{mp}}^-)^2}}_B Q_2^- \end{aligned} \quad (95)$$

where (a) follows from (90) and (87), and (b) follows from (92). From this one can obtain  $\tau_1^- = \mathbb{E}(P_1^-)^2$  which can be calculated using the knowledge that  $P_1^0, Q_2^-$  are independent Gaussian with covariances  $\mathbb{E}(P_1^0)^2 = \sigma_{\text{tr}}^2 \text{Var}(W^0)$ ,  $\mathbb{E}(Q_2^-)^2 = \sigma_d^2$ . Further,  $P_1^0, Q_2^-$  are independent of  $(S_{\text{mp}}^+, S_{\text{mp}}^-)$ .

Observe that by (95) we have

$$\tau_1^- = \frac{1}{(\bar{\alpha}_1^+)^2} \left( \mathbb{E}(A^2) \sigma_{\text{tr}}^2 \text{Var}(W^0) + \mathbb{E}(B^2) \sigma_d^2 \right). \quad (96)$$

with some simplification we get

$$\mathbb{E}(A^2) = \left( \frac{\lambda}{\sigma_{\text{tr}}^2 \beta} \right)^2 G' - \left( \frac{\lambda}{\sigma_{\text{tr}}^2 \beta} G \right)^2, \quad (97a)$$

$$\mathbb{E}(B^2) = G - \frac{\lambda}{\sigma_{\text{tr}}^2 \beta} G', \quad (97b)$$

where  $G = G_{\text{mp}}(-\frac{\lambda}{\sigma_{\text{tr}}^2 \beta})$ , with  $G_{\text{mp}}$  given in Appendix H, and  $G'$  is the derivative of  $G_{\text{mp}}$  calculated at  $-\frac{\lambda}{\sigma_{\text{tr}}^2 \beta}$ .

Now consider the **under-parametrized** case ( $\beta < 1$ ):

Let  $u = -\frac{\lambda}{\sigma_{\text{tr}}^2 \beta}$  and  $z = G_{\text{mp}}(u)$ . In this case we have

$$\bar{\alpha}_1^+ = 1 - \frac{\lambda}{\sigma_{\text{tr}}^2 \beta} G = 1 + uz. \quad (98)$$

Note that,

$$\begin{aligned} G_{\text{mp}}^{-1}(z) = u &\stackrel{(a)}{\Rightarrow} R_{\text{mp}}(z) + \frac{1}{z} = u \\ &\stackrel{(b)}{\Rightarrow} \frac{1}{1 - \beta z} + \frac{1}{z} = u, \end{aligned} \quad (99a)$$

where  $R_{\text{mp}}(\cdot)$  is the R-transform defined in (Tulino et al., 2004) and (a) follows from the relationship between the R- and Stieltjes-transform and (b) follows from the fact that for Marchenko-Pastur distribution we have  $R_{\text{mp}}(z) = \frac{1}{1 - z\beta}$ . Therefore,

$$\begin{aligned} G_{\text{mp}}\left(\frac{1}{1 - \beta z} + \frac{1}{z}\right) &= z \\ \Rightarrow G'_{\text{mp}}\left(\frac{1}{1 - \beta z} + \frac{1}{z}\right) &= G' = \frac{1}{\frac{\beta}{(1 - \beta z)^2} - \frac{1}{z^2}}. \end{aligned} \quad (100)$$

For the **over-parametrized** case ( $\beta > 1$ ) we have:

$$\bar{\alpha}_1^+ = \frac{1}{\beta} \left( 1 + \frac{\lambda}{\sigma_{\text{tr}}^2 \beta} G \right) = \frac{1 - uz}{\beta}. \quad (101)$$

In this case, as mentioned in Appendix H and following the results from (Tulino et al., 2004), the measure  $\mu_\beta$  scales with  $\beta$  and thus  $R_{\text{mp}}(z) = \frac{\beta}{1 - z}$ . Therefore, similar to (99a),  $z$  satisfies

$$\frac{\beta}{1 - z} + \frac{1}{z} = u \quad \Rightarrow G' = \frac{1}{\frac{\beta}{(1 - z)^2} - \frac{1}{z^2}}. \quad (102)$$

Now  $\tau_1^-$  can be calculated as follows:

$$\tau_1^- = \eta^2 \left( u^2 z^2 \sigma_{\text{tr}}^2 \text{var}(W^0) (\kappa - 1) + \sigma_d^2 z (uz\kappa + 1) \right) \quad (103)$$

where

$$\eta = \begin{cases} \frac{1}{(1 + uz)} & \beta < 1 \\ \frac{\beta}{(1 - uz)} & \beta \geq 1 \end{cases}, \quad \kappa = \begin{cases} \frac{(1 - \beta z)^2}{\beta z^2 - (1 - \beta z)^2} & \beta < 1 \\ \frac{(1 - z)^2}{\beta z^2 - (1 - z)^2} & \beta \geq 1 \end{cases} \quad (104)$$



and  $z$  is the solution to the fixed points

$$\begin{cases} \frac{1}{1-\beta z} + \frac{1}{z} = u & \beta < 1 \\ \frac{\beta}{1-z} + \frac{1}{z} = u & \beta \geq 1 \end{cases}. \quad (105)$$

□

### G.3. Ridgeless Linear Regression

Here we consider the case of Ridge regression (83) when  $\lambda \rightarrow 0^+$ . Note that the solution to the problem (83) is  $(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$  remains unique since  $\lambda > 0$ . The following result was stated in (Hastie et al., 2019), and can be recovered using our methodology. Note however, that we calculate the generalization error whereas they have calculated the squared error, whereby we obtain an additional additive factor of  $\sigma_d^2$ . The result explains the double-descent phenomenon for Ridgeless linear regression.

**Corollary 3.** *For ridgeless linear regression, we have*

$$\lim_{\lambda \rightarrow 0^+} \mathcal{E}_{\text{ts}}^{\text{RR}} = \begin{cases} \frac{1}{1-\beta} \sigma_d^2 & \beta < 1 \\ \frac{\beta}{\beta-1} \sigma_d^2 + (1 - \frac{1}{\beta}) \sigma_{\text{tr}}^2 \text{Var}(W^0) & \beta \geq 1 \end{cases}$$

*Proof of Corollary 3.* We calculate the parameters  $\bar{\gamma}_0^+$ ,  $\bar{\gamma}_1^-$ ,  $k_{22}$  and  $\tau_1^-$  when  $\lambda \rightarrow 0^+$ . Before starting off, we note that

$$G_0 := \lim_{z \rightarrow 0^+} G_{\text{mp}}(-z) = \begin{cases} \frac{\beta}{1-\beta} & \beta < 1 \\ \frac{\beta}{\beta-1} & \beta > 1 \end{cases}, \quad (106)$$

as described in Appendix H. Following the derivations in Corollary 2, we have

$$\bar{\gamma}_0^+ = \lambda/\beta, \quad k_{22} = \text{Var}(W^0) \quad (107)$$

Now for  $\lambda \rightarrow 0^+$ , we have

$$1 - \bar{\alpha}_1^- = \begin{cases} 1 & \beta < 1 \\ \frac{1}{\beta} & \beta \geq 1 \end{cases}, \quad \bar{\gamma}_1^- = \begin{cases} \frac{1}{G_0} = \frac{1-\beta}{\beta} & \beta < 1 \\ \frac{\lambda}{(\beta-1)\sigma_{\text{tr}}^2 \beta} & \beta > 1 \end{cases}, \quad (108)$$

Using this in simplifying (95) for  $\lambda \rightarrow 0^+$ , we get

$$\tau_1^- = \mathbb{E}(P_1^-)^2 = \begin{cases} \sigma_d^2 G_0 & \beta < 1 \\ \beta \sigma_d^2 G_0 + \sigma_{\text{tr}}^2 \text{Var}(W^0) (\beta - 1) & \beta \geq 1 \end{cases}$$

where during the evaluation of  $\mathbb{E} \left( \frac{S_{\text{mp}}^-}{\bar{\gamma}_1^- + (S_{\text{mp}}^-)^2} \right)^2$ , for the case of  $\beta > 1$ , we need to account for the point mass at 0 for  $S_{\text{mp}}^-$  with weight  $1 - \frac{1}{\beta}$ .

Next, notice that

$$a := \frac{\bar{\gamma}_0^+ \sigma_{\text{tr}}}{\bar{\gamma}_0^+ + \bar{\gamma}_1^- \sigma_{\text{tr}}^2} = \begin{cases} 0 & \beta < 1 \\ (1 - \frac{1}{\beta}) \sigma_{\text{tr}} & \beta \geq 1 \end{cases},$$

and,

$$b := \frac{\bar{\gamma}_1^- \sigma_{\text{tr}}^2}{\bar{\gamma}_0^+ + \bar{\gamma}_1^- \sigma_{\text{tr}}^2} = \begin{cases} 1 & \beta < 1 \\ \frac{1}{\beta} & \beta \geq 1 \end{cases},$$

Thus applying Corollary 1, we get

$$\begin{aligned} \mathcal{E}_{\text{ts}}^{\text{RR}} &= a^2 k_{22} + b^2 \tau_1^- + \sigma_d^2 \\ &= \begin{cases} \frac{1}{1-\beta} \sigma_d^2 & \beta < 1 \\ \frac{\beta}{\beta-1} \sigma_d^2 + (1 - \frac{1}{\beta}) \sigma_{\text{tr}}^2 \text{Var}(W^0) & \beta \geq 1 \end{cases} \end{aligned}$$

This proves the claim. □

### G.4. Train-Test Mismatch

Observe that our formulation allows for analyzing the effect of mismatch in the training and test distribution. One can consider arbitrary joint distributions over  $(S_{\text{tr}}, S_{\text{ts}})$  that model the mismatch between training and test features. Here we give a simple example which highlights the effect of this mismatch.

**Definition 4** (Bernoulli  $\varepsilon$ -mismatch).  $(S_{\text{tr}}, S_{\text{ts}})$  has a bivariate Bernoulli distribution with

- $\Pr\{S_{\text{tr}} = S_{\text{ts}} = 0\} = \Pr\{S_{\text{tr}} = S_{\text{ts}} = 1\} = (1 - \varepsilon)/2$
- $\Pr\{S_{\text{tr}} = 0, S_{\text{ts}} = 1\} = \Pr\{S_{\text{tr}} = 1, S_{\text{ts}} = 0\} = \varepsilon/2$

Notice that the marginal distribution of the  $S_{\text{tr}}$  in the Bernoulli  $\varepsilon$ -mismatch model is such that  $\Pr(S_{\text{tr}} \neq 0) = \frac{1}{2}$ . Hence half of the features extracted by the matrix  $V_0$  are relevant during training. Similarly,  $\Pr(S_{\text{ts}} \neq 0) = \frac{1}{2}$ . However the features spanned by the test data do not exactly overlap with the features captured in the training data, and the fraction of features common to both the training and test data is  $1 - \varepsilon$ . Hence for  $\varepsilon = 0$ , there is no training-test mismatch, whereas for  $\varepsilon = 1$  there is a complete mismatch.

The following result shows that the generalization error increases linearly with the mismatch parameter  $\varepsilon$ .

**Corollary 4** (Mismatch). *Consider the problem of Linear Regression (83) under the conditions of Corollary 1. Additionally suppose we have Bernoulli  $\varepsilon$ -mismatch between the training and test distributions. Then*

$$\mathcal{E}_{\text{ts}} = \frac{k_{22}}{2} ((1 - \varepsilon) \gamma^{*2} + \varepsilon) + \frac{\tau_1^-}{2} (1 - \gamma^*) (1 - \varepsilon) + \sigma_d^2,$$

where  $\gamma^* := \frac{\bar{\gamma}_0^+}{\bar{\gamma}_0^+ + \bar{\gamma}_1^-}$ . The terms  $k_{22}, \tau_1^-, \gamma^*$  are independent of  $\varepsilon$ .

*Proof.* This follows directly by calculating the expectations of the terms in Corollary 1, with the joint distribution of  $(S_{\text{tr}}, S_{\text{ts}})$  given in Definition 4. □

The quantities  $k_{22}$  and  $\tau_1^-$  in the result above can be calculated similar to the derivation in the proof of Corollary 2 and can in general depend on the regularization parameter  $\lambda$  and overparameterization parameter  $\beta$ .

### G.5. Logistic Regression

The precise analysis for the special case of regularized logistic regression estimator with i.i.d Gaussian features is provided in (Salehi et al., 2019). Consider the logistic regression model,

$$\mathbb{P}(y_i = 1 | \mathbf{x}_i) := \rho(\mathbf{x}_i^\top \mathbf{w}) \quad \text{for } i = 1, \dots, N$$

where  $\rho(x) = \frac{1}{1+e^{-x}}$  is the standard logistic function.

In this problem we consider estimates of  $\mathbf{w}^0$  such that

$$\hat{\mathbf{w}} := \underset{\mathbf{w}}{\operatorname{argmin}} \mathbf{1}^\top \log(1 + e^{\mathbf{X}\mathbf{w}}) - \mathbf{y}^\top \mathbf{X}\mathbf{w} + F_{\text{in}}(\mathbf{w}),$$

where  $F_{\text{in}}$  is the regularization function. This is a special case of optimization problem (2) where

$$F_{\text{out}}(\mathbf{y}, \mathbf{X}\mathbf{w}) = \mathbf{1}^\top \log(1 + e^{\mathbf{X}\mathbf{w}}) - \mathbf{y}^\top \mathbf{X}\mathbf{w}. \quad (109)$$

Similar to the linear regression model, using the ML-VAMP GLM learning algorithm, we can characterize the generalization error for this model with quantities  $\mathbf{K}_0^+$ ,  $\tau_1^-$ ,  $\bar{\gamma}_0^+$ ,  $\bar{\gamma}_1^-$  given by algorithm 2. We note that in this case, the output non-linearity is

$$\phi_{\text{out}}(p_2, d) = \mathbb{1}_{\{\rho(p_2) > d\}} \quad (110)$$

where  $d \sim \text{Unif}(0, 1)$ . Also, the denoisers  $g_0^+$ , and  $g_3^-$  can be derived as the proximal operators of  $F_{\text{in}}$ , and  $F_{\text{out}}$  defined in (25).

### G.6. Support Vector Machines

The asymptotic generalization error for support vector machine (SVM) is provided in (Deng et al., 2019). Our model can also handle SVMs. Similar to logistic regression, SVM finds a linear classifier using the hinge loss instead of logistic loss. Assuming the class labels are  $y = \pm 1$  the hinge loss is

$$\ell_{\text{hinge}}(y, \hat{y}) = \max(0, 1 - y\hat{y}). \quad (111)$$

Therefore, if we take

$$F_{\text{out}}(\mathbf{y}, \mathbf{X}\mathbf{w}) = \sum_i \max(0, 1 - y_i \mathbf{X}_i \mathbf{w}), \quad (112)$$

where  $\mathbf{X}_i$  is the  $i^{\text{th}}$  row of the data matrix, the ML-VAMP algorithm for GLMs finds the SVM classifier. The algorithm would have proximal map of hinge loss and our theory provides exact predictions for the estimation and prediction error of SVM.

As with all other models considered in this work, the true underlying data generating model could be anything that can be represented by the graphical model in Figure 1, e.g. logistic or probit model, and our theory is able to exactly predict the error when SVM is applied to learn such linear classifiers in the large system limit.

### H. Marchenko-Pastur distribution

We describe the random variable  $S_{\text{mp}}$  defined in (12) where  $S_{\text{mp}}^2$  has a rescaled Marchenko-Pastur distribution. Notice that the positive entries of  $\mathbf{s}_{\text{mp}}$  are the positive eigenvalues of  $\mathbf{U}^\top \mathbf{U}$  (or  $\mathbf{U}\mathbf{U}^\top$ ).

Observe that  $U_{ij} \sim N(0, \frac{1}{p})$ , whereas, the standard scaling while studying the Marchenko-Pastur distribution is for matrices  $\mathbf{H}$  such that  $H_{ij} \sim \mathcal{N}(0, \frac{1}{N})$  (for e.g. see equation (1.10) from (Tulino et al., 2004) and the discussion preceding it). Also notice that  $\sqrt{\beta}\mathbf{U}$  has the same distribution as  $\mathbf{H}$ . Thus the results from (Tulino et al., 2004) apply directly to the distributions of eigenvalues of  $\beta\mathbf{U}^\top \mathbf{U}$  and  $\beta\mathbf{U}\mathbf{U}^\top$ . We state their result below taking into account this disparity in scaling.

The positive eigenvalues of  $\beta\mathbf{U}^\top \mathbf{U}$  have an empirical distribution which converges to the following density:

$$\mu_\beta(x) = \frac{\sqrt{(b_\beta - x)_+(x - a_\beta)_+}}{2\pi\beta x} \quad (113)$$

where  $a_\beta = (1 - \sqrt{\beta})^2$ ,  $b_\beta := (1 + \sqrt{\beta})^2$ . Similarly the positive eigenvalues of  $\beta\mathbf{U}\mathbf{U}^\top$  have an empirical distribution converging to the density  $\beta\mu_\beta$ . We note the following integral which is useful in our analysis:

$$\begin{aligned} G_0 &:= \lim_{z \rightarrow 0^-} \mathbb{E} \frac{1}{S_{\text{mp}}^2 - z} \mathbb{1}_{\{S_{\text{mp}} > 0\}} \\ &= \lim_{z \rightarrow 0^-} \int_{a_\beta}^{b_\beta} \frac{1}{x/\beta - z} \mu_\beta(x) dx = \frac{\beta}{|\beta - 1|}. \end{aligned} \quad (114)$$

More generally, the Stieltjes transform of the density is given by:

$$G_{\text{mp}}(z) = \mathbb{E} \frac{1}{S_{\text{mp}}^2 - z} \mathbb{1}_{\{S_{\text{mp}} > 0\}} = \int_{a_\beta}^{b_\beta} \frac{1}{x/\beta - z} \mu_\beta(x) dx \quad (115)$$