

Functional Imaging data classification with crowdsourced noisy labels

Shuoyang WANG¹, Grace Y. YI², and Guanqun CAO³

¹ *University of Louisville, Louisville, KY 40208, USA*

² *University of Western Ontario, London, ON N6A 3K7, Canada*

³ *Michigan State University, East Lansing, MI 48824, USA*

Supplementary Material

In this supplementary material, we provide technical details, including additional assumptions, lemmas, the proofs for Theorems 1 and 2, and Proposition 1, and additional numerical results.

S1 Technical proofs

In this section, we show the proof of Theorems 1 and 2, and additional simulation results. Throughout the section, for any vectors $\mathbf{a} = (a_1, \dots, a_K)^\top$ and $\mathbf{b} = (b_1, \dots, b_K)^\top$ with length K , denote

$$\mathbf{a}^\top \log \left(\frac{\mathbf{a}}{\mathbf{b}} \right) = \sum_{k=1}^K a_k \log \left(\frac{a_k}{b_k} \right).$$

For notation simplicity and readability, we denote $\widehat{\mathbf{f}}^{\text{afDNN}} = \widehat{\mathbf{f}}^*$ as an estimator of \mathbf{f}^* .

S1.1 Additional assumptions

For function class \mathcal{H} , denote

$$(\widehat{k}, \widehat{u}) = \arg \max_{\substack{u=0, \dots, q^{(k)} \\ (k, u), k=1, \dots, K}} \frac{t_u^{(k)}}{\widetilde{\beta}_u^{(k)}}$$

as the roughest compositional function layer index pair, where

$$\widetilde{\beta}_u^{(k)} = \beta_u^{(k)} \prod_{\ell=u+1}^{q^{(k)}} \beta_\ell^{(k)} \wedge 1 \text{ and } (\widehat{l}, \widehat{v}) = \arg \max_{\substack{(k, u), k=1, \dots, K, \\ u=0, \dots, q^{(k)}}} \frac{t_u^{(k)}}{\widetilde{\beta}_u^{(k)} (1 + \alpha_k)}.$$

Let

$$\theta = \frac{(1 + \alpha_{\widehat{l}}) \widetilde{\beta}_{\widehat{v}}^{(\widehat{l})}}{(1 + \alpha_{\widehat{l}}) \widetilde{\beta}_{\widehat{v}}^{(\widehat{l})} + t_{\widehat{v}}^{(\widehat{l})}}, \quad \nu = \frac{\theta t_{\widehat{u}}^{(\widehat{k})}}{\widetilde{\beta}_{\widehat{u}}^{(\widehat{k})} (1 + \widetilde{\alpha})}, \text{ and } \widetilde{\beta} = \max_{k=1, \dots, K} (\widetilde{\beta}_0^{(k)} \wedge 1),$$

where $\widetilde{\alpha} = \min_k (\alpha_k \wedge 1)$. The following two assumptions characterize the DNN function class for densely observed functional data and discretely observed functional data, respectively, both of which are necessary to ensure fast convergence.

Assumption S1. *The DNN class \mathcal{D} satisfies*

- (a) $L \asymp \log n$;
- (b) $n^{\theta/\rho} \lesssim J \lesssim n^\nu$;
- (c) $\max_{1 \leq \ell \leq L} \gamma_\ell \asymp n^\nu$;
- (d) $s \asymp n^\nu \log n$.

Assumption S2. *Let $\underline{N} = \min_{i=1, \dots, n} N_i$ denote the least sampling frequency for the functional observations. For some phase transition point $N^* \in \mathbb{N}^+$, the bias of extracted scores satisfies $\sup_{j=1, \dots, J} E|\widehat{\xi}_{ij} - \xi_{ij}| \lesssim N_i^{-\tau}$, and the DNN class \mathcal{D} satisfies*

$$(a) L \asymp (\log n)\mathbb{I}(\underline{N} \geq N^*) + (\log \underline{N})\mathbb{I}(\underline{N} < N^*);$$

$$(b) n^{\theta/\rho}\mathbb{I}(\underline{N} \geq N^*) + \underline{N}^{\theta'/\rho}\mathbb{I}(\underline{N} < N^*) \lesssim J \lesssim n^\nu\mathbb{I}(\underline{N} \geq N^*) + \underline{N}^{\nu'}\mathbb{I}(\underline{N} < N^*);$$

$$(c) \max_{1 \leq \ell \leq L} \gamma_\ell \asymp n^\nu\mathbb{I}(\underline{N} \geq N^*) + \underline{N}^{\nu'}\mathbb{I}(\underline{N} < N^*);$$

$$(d) s \asymp (n^\nu \log n)\mathbb{I}(\underline{N} \geq N^*) + (\underline{N}^{\nu'} \log \underline{N})\mathbb{I}(\underline{N} < N^*).$$

$$\text{where } \theta' = \tau \tilde{\beta}, \text{ and } \nu' = \frac{\theta'_t(\hat{k})}{\tilde{\beta}_u^{(k)}(1+\tilde{\alpha})}.$$

S1.2 Technical lemmas

For any fixed J , define the training loss difference between $\hat{\mathbf{f}}^*$ and the global minimum over \mathcal{D} with respect to the true labels Y_i as

$$\delta_n(J) = \mathbb{E} \left\{ -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K \mathbb{I}(Y_i = k) \log \hat{f}_k^*(\boldsymbol{\xi}_{iJ}) - \min_{\mathbf{f} \in \mathcal{D}} -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K \mathbb{I}(Y_i = k) \log f_k(\boldsymbol{\xi}_{iJ}) \right\}.$$

Lemma S1. *If a network $\hat{\mathbf{f}}^* \in \mathcal{D}(L, J, \mathbf{p}, s)$ and $\mathcal{D}(L, J, \mathbf{p}, s)$ satisfies*

- $L \asymp \log n$,
- $\|\mathbf{p}\|_{\max} \asymp \max\{J, K, C_2 n^\nu\}$,
- $s \asymp C_3 n^\nu \log n$,

then

$$\mathbb{E}_{\mathbf{f}^*} \left\{ \mathbf{f}^{*\top}(\boldsymbol{\xi}_J) \left[(C_0 \epsilon^{-1}) \wedge \log \left(\frac{\mathbf{f}^*(\boldsymbol{\xi}_J)}{\hat{\mathbf{f}}^*(\boldsymbol{\xi}_J)} \right) \right] \right\} \lesssim n^{-\theta} \log^3 n$$

for $\hat{\mathbf{f}}^*$ taking values on $[a, b]$, $a < b \in \mathbb{R}$, and ϵ is defined in Assumption 2(a).

Lemma S2. (*Bos and Schmidt-Hieber (2022) Theorem 3.5*) *Given a fixed $J \in \mathbb{N}^+$, for*

any classifier $\hat{\mathbf{f}}^*$, we have

$$\begin{aligned} & \mathbb{E}_{\mathbf{f}^*} \left\{ \mathbf{f}^{*\top}(\boldsymbol{\xi}_J) \left[C_0 \wedge \log \left(\frac{\mathbf{f}^*(\boldsymbol{\xi}_J)}{\hat{\mathbf{f}}^*(\boldsymbol{\xi}_J)} \right) \right] \right\} \\ & \leq C^* n^{-\theta} \log^3 n + \delta_n(J) \\ & \quad + 4n^{-1} \left[68C_0(s+1) \log \left(2^{2L+6} n(L+1) K^3 J^2 s^L \right) + 272C_0 + 2K(\log n + C_0) \right], \end{aligned}$$

where C^* is some universal constant.

In the sequel, define the empirical risk with respect to the proposed loss function (3.4):

$$\begin{aligned} \tilde{\delta}_n(J) &= \mathbb{E} \left\{ -\frac{1}{nR} \sum_{i=1}^n \sum_{r=1}^R \sum_{k=1}^K \mathbb{I}(\tilde{Y}_i = k) \log \left[\hat{\mathbf{A}}^{*(r)} \hat{\mathbf{f}}^*(\boldsymbol{\xi}_{iJ}) \right]_k + \lambda \operatorname{tr}(\hat{\mathbf{A}}^{*(r)}) \right. \\ & \quad \left. - \min_{\mathbf{f} \in \mathcal{D}, \mathbf{A}^{(r)} \in \mathcal{A}} -\frac{1}{nR} \sum_{r=1}^R \sum_{i=1}^n \sum_{k=1}^K \mathbb{I}(\tilde{Y}_i^{(r)} = k) \log[\mathbf{A}^{(r)} \mathbf{f}(\boldsymbol{\xi}_{iJ})]_k + \lambda \operatorname{tr}(\mathbf{A}^{(r)}) \right\}, \end{aligned}$$

where J is considered as a fixed positive integer.

The following lemma shows an information upper bound of the KL divergence with respect to the classifier derived from the proposed loss function (3.4).

Lemma S3. *Assume Assumption 1. For any $J \in \mathbb{N}^+$, and for any classifier $\hat{\mathbf{f}}^*$ obtained by the loss function in Equation (3.4), when $\boldsymbol{\xi}_J$ is defined on some $[-M, M]^J \subset \mathbb{R}$ for $M > 0$, we have*

$$\begin{aligned} & \mathbb{E} \left\{ \mathbf{f}^{*\top}(\boldsymbol{\xi}_J) \left[(C_0 \epsilon^{-1}) \wedge \log \left(\frac{\mathbf{f}^*(\boldsymbol{\xi}_J)}{\hat{\mathbf{f}}^*(\boldsymbol{\xi}_J)} \right) \right] \right\} \\ & \lesssim n^{-\theta} \log^3 n + \tilde{\delta}_n(J) + R^{-1} \sum_{r=1}^R \operatorname{tr}(\mathbf{A}^{*(r)} - \mathbf{I}_K) \\ & \quad + n^{-1} \left[s \log \left(2^{2L+6} n(L+1) K^3 J^2 s^L \right) + K \log n \right]. \end{aligned}$$

Proof. First, for any $\mathbf{f} = (f_k)_{k=1,\dots,K}$ and $\mathbf{A}^{(r)}, r = 1, \dots, R$, it satisfies that

$$\begin{aligned} \mathbb{E} \left(\sum_{k=1}^K \mathbb{I}(Y_i = k) \log f_k(\boldsymbol{\xi}_{iJ}) \right) &= \mathbb{E} \left(\sum_{k=1}^K \mathbb{E}(\mathbb{I}(Y_i = k) | \boldsymbol{\xi}_{iJ}) \log f_k(\boldsymbol{\xi}_{iJ}) \right) \\ &= \mathbb{E} \left(\sum_{k=1}^K f_k^*(\boldsymbol{\xi}_{iJ}) \log f_k(\boldsymbol{\xi}_{iJ}) \right), \end{aligned}$$

and

$$\begin{aligned} &\mathbb{E} \left(\sum_{k=1}^K \mathbb{I}(\tilde{Y}_i^{(r)} = k) \log [\mathbf{A}^{(r)} \mathbf{f}(\boldsymbol{\xi}_{iJ})]_k \right) \\ &= \mathbb{E} \left(\sum_{k=1}^K \mathbb{E}(\mathbb{I}(\tilde{Y}_i^{(r)} = k) | \boldsymbol{\xi}_{iJ}) \log [\mathbf{A}^{(r)} \mathbf{f}(\boldsymbol{\xi}_{iJ})]_k \right) \\ &= \mathbb{E} \left(\sum_{k=1}^K [\mathbf{A}^{*(r)} \mathbf{f}^*(\boldsymbol{\xi}_{iJ})]_k \log [\mathbf{A}^{(r)} \mathbf{f}(\boldsymbol{\xi}_{iJ})]_k \right). \end{aligned}$$

Next,

$$\begin{aligned} &(\mathbf{A}^{*(r)} \mathbf{f}^*(\boldsymbol{\xi}_{iJ}))^\top \log(\mathbf{A}^{(r)} \mathbf{f}(\boldsymbol{\xi}_{iJ})) - (\mathbf{f}^*(\boldsymbol{\xi}_{iJ}))^\top \log(\mathbf{f}(\boldsymbol{\xi}_{iJ})) \\ &= (\mathbf{f}^*)^\top \left[(\mathbf{A}^{*(r)})^\top \log(\mathbf{A}^{(r)} \mathbf{f}(\boldsymbol{\xi}_{iJ})) - \mathbf{I}_k \log(\mathbf{f}(\boldsymbol{\xi}_{iJ})) \right] \\ &= (\mathbf{f}^*)^\top \left[(\mathbf{A}^{*(r)})^\top \left\{ \log(\mathbf{f}(\boldsymbol{\xi}_{iJ})) + \frac{\mathbf{A}^{(r)} \mathbf{f}(\boldsymbol{\xi}_{iJ}) - \mathbf{f}(\boldsymbol{\xi}_{iJ})}{\mathbf{A}^{(r)} \mathbf{f}(\boldsymbol{\xi}_{iJ})} + o(\mathbf{A}^{(r)} \mathbf{f}(\boldsymbol{\xi}_{iJ}) - \mathbf{f}(\boldsymbol{\xi}_{iJ})) \right\} \right. \\ &\quad \left. - \mathbf{I}_k \log(\mathbf{f}(\boldsymbol{\xi}_{iJ})) \right] \\ &= (\mathbf{f}^*)^\top \left[(\mathbf{A}^{*(r)} - \mathbf{I}_k)^\top \log(\mathbf{f}(\boldsymbol{\xi}_{iJ})) + (\mathbf{A}^{*(r)})^\top \frac{\mathbf{A}^{(r)} \mathbf{f}(\boldsymbol{\xi}_{iJ}) - \mathbf{f}(\boldsymbol{\xi}_{iJ})}{\mathbf{A}^{(r)} \mathbf{f}(\boldsymbol{\xi}_{iJ})} \right. \\ &\quad \left. + o(\mathbf{A}^{(r)} \mathbf{f}(\boldsymbol{\xi}_{iJ}) - \mathbf{f}(\boldsymbol{\xi}_{iJ})) \right], \end{aligned}$$

where the second equality is derived by the Taylor series expansion. Therefore, under

Assumption 1 and Proposition 1, when $\mathbf{A}^{(r)} = \widehat{\mathbf{A}}^{*(r)}$ and $\mathbf{f} = \widehat{\mathbf{f}}^*$, since $\|\mathbf{f}^*\|_{\max} \leq 1$

and $\|\mathbf{f}\|_{\max} \leq 1$, it follows that

$$(\mathbf{A}^{*(r)} \mathbf{f}^*(\boldsymbol{\xi}_{iJ}))^\top \log(\mathbf{A}^{(r)} \mathbf{f}(\boldsymbol{\xi}_{iJ})) - (\mathbf{f}^*(\boldsymbol{\xi}_{iJ}))^\top \log(\mathbf{f}(\boldsymbol{\xi}_{iJ})) = O(\text{tr}(\mathbf{A}^{*(r)} - \mathbf{I}_K)).$$

The same bound can also be derived by letting $\mathbf{A}^{(r)} = \widehat{\mathbf{A}}^{*(r)}$ and

$$\mathbf{f} = \arg \min_{\mathbf{f}' \in \mathcal{D}} -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K \mathbb{I}(Y_i = k) \log f'_k(\boldsymbol{\xi}_{iJ}).$$

As a result, we obtain that

$$\widetilde{\delta}(J) - \delta(J) = O\left(R^{-1} \sum_{r=1}^R \text{tr}(\mathbf{A}^{*(r)} - \mathbf{I}_K)\right),$$

where the result holds under Assumption 1. Substituting the $\delta(J)$ by $\widetilde{\delta}(J) - \widetilde{\delta}(J) + \delta(J)$ in Lemma S2 and adjusting the truncation constant as $C_0\epsilon^{-1}$, we complete the proof, where the constant ϵ is defined in Assumption 2. \square

S1.3 Proof of Proposition 1

Proof. According to Theorem 1 in Tanno et al. (2019), it suffices to show that for any positive λ ,

$$(\mathbf{A}^{*(1)}, \dots, \mathbf{A}^{*(R)}) = \arg \min_{\{\widehat{\mathbf{A}}^{*(r)} \in \mathcal{A}\}_{r=1}^R} \lambda \text{tr} \left(\sum_{r=1}^R \widehat{\mathbf{A}}^{*(r)} \right)$$

and this solution is unique. Since $(\mathbf{A}^{*(1)}, \dots, \mathbf{A}^{*(R)})$ is the maximizer of the likelihood function, $(\mathbf{A}^{*(1)}, \dots, \mathbf{A}^{*(R)})$ is the minimizer of the objective function. The uniqueness is straightforward, as if the solution to

$$\arg \min_{\{\widehat{\mathbf{A}}^{*(r)} \in \mathcal{A}\}_{r=1}^R} -\frac{1}{nR} \sum_{r=1}^R \sum_{i=1}^n \sum_{k=1}^K \mathbb{I}(\widetilde{Y}_i^{(r)} = k) \log[\widehat{\mathbf{A}}^{*(r)} \mathbf{f}^*(\boldsymbol{\xi}_J^{(i)})]_k$$

is not unique, the uniqueness of $\arg \min_{(\widehat{\mathbf{A}}^{*(1)}, \dots, \widehat{\mathbf{A}}^{*(R)})} \text{tr} \left(\sum_{r=1}^R \widehat{\mathbf{A}}^{*(r)} \right)$ still ensures the solutions are $(\mathbf{A}^{*(1)}, \dots, \mathbf{A}^{*(R)})$, thus the proof is complete. \square

S1.4 Proof of Theorem 1

Let the set of the effective inputs of the k th be grouped by

$$\mathcal{P}_k = \left\{ j : \xi_j \text{ is effective for } g_{0v}^{(k)} \text{ for all } v = 1, \dots, d_1^{(k)} \right\}.$$

By definition of \mathcal{G} , $|\mathcal{A}_k| \leq t_0^{(k)} d_1^{(k)} < \infty$. Let $\mathcal{P} = \bigcup_{k=1}^K \mathcal{P}_k$ denote the index set of effective inputs among all K groups, where $|\mathcal{P}| \leq \sum_{k=1}^K t_0^{(k)} d_1^{(k)} < \infty$. Given any positive integer m , define $\mathcal{Q}_M = \{|\xi_j| \leq M \text{ for all } j \in \mathcal{P}\}$ to be the concentration set for all effective inputs.

We first focus on the convergence rate by considering \mathcal{Q}_M and \mathcal{Q}_M^c separately. Let M be relatively large, and write identity vector $\mathbf{1}_K = (1, \dots, 1)^\top$ with length K . For any J , we have

$$\begin{aligned} & \mathbb{E} \left[\mathbf{f}^{\otimes}(\boldsymbol{\xi})^\top \left\{ C_0 \mathbf{1}_K \wedge \log \left(\frac{\mathbf{f}^{\otimes}(\boldsymbol{\xi})}{\widehat{\mathbf{f}}^*(\boldsymbol{\xi}_J)} \right) \right\} \right] \\ &= \mathbb{E} \left[\left(\mathbf{f}^{\otimes}(\boldsymbol{\xi})^\top C_0 \right) \wedge \left\{ \mathbf{f}^{\otimes}(\boldsymbol{\xi})^\top \log \left(\frac{\mathbf{f}^{\otimes}(\boldsymbol{\xi})}{\widehat{\mathbf{f}}^*(\boldsymbol{\xi}_J)} \right) \right\} \right] \\ &\leq \mathbb{E} \left\{ C_0 \wedge \left[\left(\mathbf{f}^{\otimes}(\boldsymbol{\xi}) - \mathbf{f}^*(\boldsymbol{\xi}_J) + \mathbf{f}^*(\boldsymbol{\xi}_J) \right)^\top \log \left(\frac{\mathbf{f}^{\otimes}(\boldsymbol{\xi})}{\mathbf{f}^*(\boldsymbol{\xi}_J)} \cdot \frac{\mathbf{f}^*(\boldsymbol{\xi}_J)}{\widehat{\mathbf{f}}^*(\boldsymbol{\xi}_J)} \right) \right] \right\} \\ &\leq \mathbb{E} \left[\mathbf{f}^{\otimes}(\boldsymbol{\xi})^\top \log \left(\frac{\mathbf{f}^{\otimes}(\boldsymbol{\xi})}{\mathbf{f}^*(\boldsymbol{\xi}_J)} \right) \right] + \mathbb{E} \left[C_0 \wedge \left| \left(\mathbf{f}^{\otimes}(\boldsymbol{\xi}) - \mathbf{f}^*(\boldsymbol{\xi}_J) \right)^\top \log \left(\frac{\mathbf{f}^*(\boldsymbol{\xi}_J)}{\widehat{\mathbf{f}}^*(\boldsymbol{\xi}_J)} \right) \right| \right] \\ &\quad + \mathbb{E} \left\{ C_0 \wedge \left[\mathbf{f}^{*\top}(\boldsymbol{\xi}_J) \log \left(\frac{\mathbf{f}^*(\boldsymbol{\xi}_J)}{\widehat{\mathbf{f}}^*(\boldsymbol{\xi}_J)} \right) \right] \right\} \\ &\leq \sum_{k=1}^K \mathbb{E} \left[\frac{(f_k^{\otimes}(\boldsymbol{\xi}) - f_k^*(\boldsymbol{\xi}_J))^2}{f_k^*(\boldsymbol{\xi}_J)} \right] \\ &\quad + \sum_{k=1}^K \mathbb{E} \left[C_0 \wedge \left| (f_k^{\otimes}(\boldsymbol{\xi}) - f_k^*(\boldsymbol{\xi}_J)) \log \left(1 + \frac{f_k^*(\boldsymbol{\xi}_J) - \widehat{f}^*(\boldsymbol{\xi}_J)}{\widehat{f}^*(\boldsymbol{\xi}_J)} \right) \right| \mathbb{I}(\mathcal{Q}_M) \right] \\ &\quad + \mathbb{E} \left\{ C_0 \wedge \left[\mathbf{f}^{*\top}(\boldsymbol{\xi}_J) \log \left(\frac{\mathbf{f}^*(\boldsymbol{\xi}_J)}{\widehat{\mathbf{f}}^*(\boldsymbol{\xi}_J)} \right) \right] \mathbb{I}(\mathcal{Q}_M) \right\} + 2C_0 K \mathbb{P}(\mathcal{Q}_M^c), \tag{S1} \end{aligned}$$

where the first inequality holds because of $f_k^{\otimes} \leq 1$, the second inequality holds since for any $a, b, c \geq 0$, we have that $a \wedge (b+c) \leq (a+c) \wedge (b+c) = (a \wedge b) + c$. The last inequality holds for the first term since KL-divergence is upper bounded by χ^2 -divergence.

For the first term in (S1), when J is relatively large, for defined in Assumption 2, we have

$$\begin{aligned}
& \sum_{k=1}^K \mathbb{E} \left[\frac{(f_k^{\otimes}(\boldsymbol{\xi}) - f_k^*(\boldsymbol{\xi}_J))^2}{f_k^*(\boldsymbol{\xi}_J)} \right] \\
&= \sum_{k=1}^K \mathbb{E} \left[\frac{(f_k^{\otimes}(\boldsymbol{\xi}) - f_k^*(\boldsymbol{\xi}_J))^2}{f_k^*(\boldsymbol{\xi}_J)} \mathbb{I}(f_k^*(\boldsymbol{\xi}_J) > \epsilon) \right] \\
&\quad + \sum_{k=1}^K \mathbb{E} \left[\frac{(f_k^{\otimes}(\boldsymbol{\xi}) - f_k^*(\boldsymbol{\xi}_J))^2}{f_k^*(\boldsymbol{\xi}_J)} \mathbb{I}(f_k^*(\boldsymbol{\xi}_J) \leq \epsilon) \right] \\
&< \sum_{k=1}^K \mathbb{E} \left[\frac{(f_k^{\otimes}(\boldsymbol{\xi}) - f_k^*(\boldsymbol{\xi}_J))^2}{\epsilon} \right] + \sum_{k=1}^K \mathbb{E} \left[\frac{(f_k^{\otimes}(\boldsymbol{\xi}) - f_k^*(\boldsymbol{\xi}_J))^2}{f_k^*(\boldsymbol{\xi}_J)} \right] \mathbb{P}(f_k^*(\boldsymbol{\xi}_J) \leq \epsilon) \\
&\lesssim \zeta(J),
\end{aligned}$$

where the second asymptotic inequality is derived by using Lemmas 7.7 and 7.8 in Wang and Cao (2024) and Theorem 3.2 in Bos and Schmidt-Hieber (2022).

For any $a, b \in \mathbb{R}$ and fixed C_0 , there exists C'_0 , such that

$$|a \log(1 + b)| \wedge C_0 \leq |a \log(1 + b)| \mathbb{I}(|b| \leq C'_0).$$

Therefore, by the fact that given M , within the concentrated set \mathcal{Q}_M , there exists an M^* such that $\left[\widehat{f}^{(k)}(\boldsymbol{\xi}_J) \right]^{-1} \mathbb{I}(\mathcal{Q}_M) \leq M^*$ for any $\boldsymbol{\xi}_J$, the second term in (S1) can be

upper bounded as

$$\begin{aligned}
& \sum_{k=1}^K \mathbb{E} \left[C_0 \wedge \left| (f_k^\otimes(\boldsymbol{\xi}) - f_k^*(\boldsymbol{\xi}_J)) \log \left(1 + \frac{f_k^*(\boldsymbol{\xi}_J) - \widehat{f}_k(\boldsymbol{\xi}_J)}{\widehat{f}_k(\boldsymbol{\xi}_J)} \right) \right| \mathbb{I}(\mathcal{Q}_M) \right] \\
& \leq \sum_{k=1}^K \mathbb{E} \left[|f_k^\otimes(\boldsymbol{\xi}) - f_k^*(\boldsymbol{\xi}_J)| \log \left(1 + \frac{|f_k^*(\boldsymbol{\xi}_J) - \widehat{f}_k(\boldsymbol{\xi}_J)| \wedge C'_0}{\widehat{f}_k(\boldsymbol{\xi}_J)} \right) \mathbb{I}(\mathcal{Q}_M) \right] \\
& \leq \sum_{k=1}^K \mathbb{E} \left[|f_k^\otimes(\boldsymbol{\xi}) - f_k^*(\boldsymbol{\xi}_J)| \frac{|f_k^*(\boldsymbol{\xi}_J) - \widehat{f}_{(k)}(\boldsymbol{\xi}_J)| \wedge C'_0}{\widehat{f}_{(k)}(\boldsymbol{\xi}_J)} \mathbb{I}(\mathcal{Q}_M) \right] \\
& \lesssim \sum_{k=1}^K \mathbb{E} (|f_k^\otimes(\boldsymbol{\xi}) - f_k^*(\boldsymbol{\xi}_J)| M^*) \\
& \lesssim \zeta(J),
\end{aligned}$$

where the second last inequality comes from Lemmas 7.5 and 7.6 in [Wang and Cao \(2024\)](#).

For any $\delta > 0$, there exists an $M = m > 0$ such that

$$\mathbb{P}(\mathcal{B}_m) = \mathbb{P}(|\xi_j| \leq m \text{ for all } j \in \mathcal{P}) \geq 1 - \delta, \text{ for all } j \in \mathcal{P}.$$

Therefore, for any $J \geq J_0$, where J_0 is defined in Assumption 3, when δ is relatively small, there exists a corresponding m , such that

$$\begin{aligned}
& \mathbb{E} \left[\mathbf{f}^\otimes(\boldsymbol{\xi})^\top \left(C_0 \wedge \log \left(\frac{\mathbf{f}^\otimes(\boldsymbol{\xi})}{\widehat{\mathbf{f}}^*(\boldsymbol{\xi}_J)} \right) \right) \right] \\
& \lesssim \zeta(J) + \mathbb{E} \left\{ \mathbf{f}^{*\top}(\boldsymbol{\xi}_J) \left[(C_0 \epsilon^{-1}) \wedge \log \left(\frac{\mathbf{f}^*(\boldsymbol{\xi}_J)}{\widehat{\mathbf{f}}^*(\boldsymbol{\xi}_J)} \right) \right] \mathbb{I}(\mathcal{Q}_m) \right\}.
\end{aligned}$$

By Lemma S3, we obtain that

$$\begin{aligned} & \mathbb{E} \left[\mathbf{f}^{\otimes}(\boldsymbol{\xi})^\top \left(C_0 \wedge \log \left(\frac{\mathbf{f}^{\otimes}(\boldsymbol{\xi})}{\widehat{\mathbf{f}}^*(\boldsymbol{\xi}_J)} \right) \right) \right] \\ & \lesssim n^{-\theta} \log^3 n + \zeta(J) + \widetilde{\delta}_n(J) + R^{-1} \sum_{r=1}^R \text{tr}(\mathbf{A}^{*(r)} - \mathbf{I}_K) \\ & \quad + n^{-1} [s \log(2^{2L+6} n(L+1) K^3 J^2 s^L) + K \log n]. \end{aligned}$$

When Assumption S1 holds, the inequality can be simplified as

$$\begin{aligned} & \mathbb{E} \left[\mathbf{f}^{\otimes}(\boldsymbol{\xi})^\top \left(C_0 \wedge \log \left(\frac{\mathbf{f}^{\otimes}(\boldsymbol{\xi})}{\widehat{\mathbf{f}}^*(\boldsymbol{\xi}_J)} \right) \right) \right] \\ & \lesssim n^{-\theta} \log^3 n + \widetilde{\delta}_n(J) + R^{-1} \sum_{r=1}^R \text{tr}(\mathbf{A}^{*(r)} - \mathbf{I}_K), \end{aligned}$$

where the $\widetilde{\delta}_n(J)$ can be sufficiently small when the loss function (3.4) is minimized well.

Thus, the proof is completed.

S1.5 Proof of Theorem 2

We first identify an the upper bound of $\mathbb{P}\left(f_k^*(\widehat{\boldsymbol{\xi}}_J) \leq x\right)$. By the definition of f_k^* , we have

$$\left| f_k^*(\boldsymbol{\xi}_{iJ}) - f_k^*(\widehat{\boldsymbol{\xi}}_{iJ}) \right| \lesssim \|\boldsymbol{\xi}_{iJ} - \widehat{\boldsymbol{\xi}}_{iJ}\|_{\max}^{\widetilde{\beta}_0^{(k)} \wedge 1},$$

and thus

$$\mathbb{E} \left| f_k^*(\boldsymbol{\xi}_{iJ}) - f_k^*(\widehat{\boldsymbol{\xi}}_{iJ}) \right| \lesssim N_i^{-\tau \widetilde{\beta}}, \text{ for any } k = 1, \dots, K.$$

Therefore, it suffices to show that

$$\mathbb{E} \left| f_k^*(\boldsymbol{\xi}_J) - f_k^*(\widehat{\boldsymbol{\xi}}_J) \right| \lesssim \underline{N}^{-\tau \widetilde{\beta}}, \forall k = 1, \dots, K.$$

According to Lemma 7.7 in Wang and Cao (2024), we have

$$\left| E \left(f_k^\circledast(\boldsymbol{\xi}) - f_k^*(\widehat{\boldsymbol{\xi}}_J) \right) \right| \lesssim \zeta(J) + \underline{N}^{-\tau\tilde{\beta}}, \text{ for any } k = 1, \dots, K.$$

According to Assumptions 2 and 3, we have for any $k = 1, \dots, K$,

$$\begin{aligned} & \mathbb{P} \left(f_k^*(\widehat{\boldsymbol{\xi}}_J) \leq x \right) \\ &= \mathbb{P} \left(f_k^*(\boldsymbol{\xi}_J) \leq x + f_k^*(\boldsymbol{\xi}_J) - f_k^*(\widehat{\boldsymbol{\xi}}_J) \right) \\ &\leq \mathbb{P} \left(f_k^*(\boldsymbol{\xi}_J) \leq x + \left| f_k^*(\boldsymbol{\xi}_J) - f_k^*(\widehat{\boldsymbol{\xi}}_J) \right| \right) \\ &\leq \mathbb{P} \left(f_k^*(\boldsymbol{\xi}_J) \leq x + \left| f_k^*(\boldsymbol{\xi}_J) - f_k^*(\widehat{\boldsymbol{\xi}}_J) \right|, \left| f_k^*(\boldsymbol{\xi}_J) - f_k^*(\widehat{\boldsymbol{\xi}}_J) \right| \leq x \right) \\ &\quad + \mathbb{P} \left(f_k^*(\boldsymbol{\xi}_J) \leq x + \left| f_k^*(\boldsymbol{\xi}_J) - f_k^*(\widehat{\boldsymbol{\xi}}_J) \right|, \left| f_k^*(\boldsymbol{\xi}_J) - f_k^*(\widehat{\boldsymbol{\xi}}_J) \right| > x \right) \\ &\leq \mathbb{P} \left(f_k^*(\boldsymbol{\xi}_J) \leq 2x \right) + \mathbb{P} \left(\left| f_k^*(\boldsymbol{\xi}_J) - f_k^*(\widehat{\boldsymbol{\xi}}_J) \right| > x \right) \\ &\lesssim x^{\alpha_k} + \zeta(J) + x \mathbb{E} \left| f_k^*(\boldsymbol{\xi}_J) - f_k^*(\widehat{\boldsymbol{\xi}}_J) \right| \\ &\lesssim x^{\alpha_k} + \zeta(J) + x\zeta(J) + x\underline{N}^{-\tau\tilde{\beta}} \\ &\lesssim x^{\alpha_k} + \zeta(J) + x\underline{N}^{-\tau\tilde{\beta}}. \end{aligned}$$

Next, we adapt the proof of Theorem 1 to bound the decomposed risk:

$$\begin{aligned} & \mathbb{E} \left[\mathbf{f}^\circledast(\boldsymbol{\xi})^\top \left\{ C_0 \mathbf{1}_K \wedge \log \left(\frac{\mathbf{f}^\circledast(\boldsymbol{\xi})}{\widehat{\mathbf{f}}^*(\widehat{\boldsymbol{\xi}}_J)} \right) \right\} \right] \tag{S2} \\ &\leq \sum_{k=1}^K \mathbb{E} \left[\frac{(f_k^\circledast(\boldsymbol{\xi}) - f_k^*(\boldsymbol{\xi}_J))^2}{f_k^*(\widehat{\boldsymbol{\xi}}_J)} \right] \\ &\quad + \sum_{k=1}^K \mathbb{E} \left[C_0 \wedge \left| (f_k^\circledast(\boldsymbol{\xi}) - f_k^*(\widehat{\boldsymbol{\xi}}_J)) \log \left(1 + \frac{f_k^*(\widehat{\boldsymbol{\xi}}_J) - \widehat{f}_k^*(\widehat{\boldsymbol{\xi}}_J)}{\widehat{f}_k^*(\widehat{\boldsymbol{\xi}}_J)} \right) \right| \mathbb{I}(\mathcal{Q}_M) \right] \\ &\quad + \mathbb{E} \left\{ C_0 \wedge \left[\mathbf{f}^{*\top}(\widehat{\boldsymbol{\xi}}_J) \log \left(\frac{\mathbf{f}^*(\widehat{\boldsymbol{\xi}}_J)}{\widehat{\mathbf{f}}^*(\widehat{\boldsymbol{\xi}}_J)} \right) \mathbb{I}(\mathcal{Q}_M) \right] \right\} + 2C_0 K \mathbb{P}(\mathcal{Q}_M^c). \end{aligned}$$

By the upper bound of $\mathbb{P} \left(f_k^*(\widehat{\boldsymbol{\xi}}_J) \right)$ in (S2), when J is relatively large, we can bound

the first term in (S2). When J is relatively large, we have

$$\begin{aligned}
& \sum_{k=1}^K \mathbb{E} \left[\frac{\left(f_k^{\otimes}(\boldsymbol{\xi}) - f_k^*(\widehat{\boldsymbol{\xi}}_J) \right)^2}{f_k^*(\widehat{\boldsymbol{\xi}}_J)} \right] \\
&= \sum_{k=1}^K \mathbb{E} \left[\frac{\left(f_k^{\otimes}(\boldsymbol{\xi}) - f_k^*(\widehat{\boldsymbol{\xi}}_J) \right)^2}{f_k^*(\widehat{\boldsymbol{\xi}}_J)} \mathbb{I} \left(f_k^*(\widehat{\boldsymbol{\xi}}_J) > \epsilon \right) \right] \\
&\quad + \sum_{k=1}^K \mathbb{E} \left[\frac{\left(f_k^{\otimes}(\boldsymbol{\xi}) - f_k^*(\widehat{\boldsymbol{\xi}}_J) \right)^2}{f_k^*(\widehat{\boldsymbol{\xi}}_J)} \mathbb{I} \left(f_k^*(\widehat{\boldsymbol{\xi}}_J) \leq \epsilon \right) \right] \\
&< \sum_{k=1}^K \mathbb{E} \left[\frac{\left(f_k^{\otimes}(\boldsymbol{\xi}) - f_k^*(\widehat{\boldsymbol{\xi}}_J) \right)^2}{\epsilon} \right] + \sum_{k=1}^K \mathbb{E} \left[\frac{\left(f_k^{\otimes}(\boldsymbol{\xi}) - f_k^*(\widehat{\boldsymbol{\xi}}_J) \right)^2}{f_k^*(\widehat{\boldsymbol{\xi}}_J)} \right] \mathbb{P} \left(f_k^*(\widehat{\boldsymbol{\xi}}_J) \leq \epsilon \right) \\
&\lesssim \zeta(J) + \underline{N}^{-\tau\tilde{\beta}}.
\end{aligned}$$

Similarly as indicated in proof of Theorem 1, by the fact that given M , within the concentrated set \mathcal{Q}_M , there exists an M^{**} such that

$$\left[\widehat{f}^{(k)}(\boldsymbol{\xi}_J) \right]^{-1} \mathbb{I}(\mathcal{Q}_M) \leq M^{**}$$

for any $\boldsymbol{\xi}_J$, the second term in (S2) can be upper bounded:

$$\begin{aligned}
& \sum_{k=1}^K \mathbb{E} \left[C_0 \wedge \left| \left(f_k^{\otimes}(\boldsymbol{\xi}) - f_k^*(\widehat{\boldsymbol{\xi}}_J) \right) \log \left(1 + \frac{f_k^*(\widehat{\boldsymbol{\xi}}_J) - \widehat{f}_k(\widehat{\boldsymbol{\xi}}_J)}{\widehat{f}_k(\widehat{\boldsymbol{\xi}}_J)} \right) \right| \mathbb{I}(\mathcal{Q}_M) \right] \\
&\leq \sum_{k=1}^K \mathbb{E} \left[\left| f_k^{\otimes}(\boldsymbol{\xi}) - f_k^*(\widehat{\boldsymbol{\xi}}_J) \right| \log \left(1 + \frac{\left| f_k^*(\widehat{\boldsymbol{\xi}}_J) - \widehat{f}_k(\widehat{\boldsymbol{\xi}}_J) \right| \wedge C'_0}{\widehat{f}_k(\widehat{\boldsymbol{\xi}}_J)} \right) \mathbb{I}(\mathcal{Q}_M) \right] \\
&\leq \sum_{k=1}^K \mathbb{E} \left[\left| f_k^{\otimes}(\boldsymbol{\xi}) - f_k^*(\widehat{\boldsymbol{\xi}}_J) \right| \frac{\left| f_k^*(\widehat{\boldsymbol{\xi}}_J) - \widehat{f}_k(\widehat{\boldsymbol{\xi}}_J) \right| \wedge C'_0}{\widehat{f}_k(\widehat{\boldsymbol{\xi}}_J)} \mathbb{I}(\mathcal{Q}_M) \right] \\
&\lesssim \sum_{k=1}^K \mathbb{E} \left(\left| f_k^{\otimes}(\boldsymbol{\xi}) - f_k^*(\widehat{\boldsymbol{\xi}}_J) \right| M^{**} \right) \\
&\lesssim \zeta(J) + \underline{N}^{-\tau\tilde{\beta}}.
\end{aligned}$$

Consequently, under Assumption S2, we obtain that

$$\begin{aligned} & \mathbb{E} \left[\mathbf{f}^{\otimes}(\boldsymbol{\xi})^\top \left(C_0 \wedge \log \left(\frac{\mathbf{f}^{\otimes}(\boldsymbol{\xi})}{\widehat{\mathbf{f}}^*(\widehat{\boldsymbol{\xi}}_J)} \right) \right) \right] \\ & \lesssim n^{-\theta} \log^3 n + \underline{N}^{-\tau\tilde{\beta}} + \tilde{\delta}_n(J) + R^{-1} \sum_{r=1}^R \text{tr}(\mathbf{A}^{*(r)} - \mathbf{I}_K), \end{aligned}$$

where the phase transition point is at $N^* \asymp (n^\theta / \log^3 n)^{1/\theta'}$, and $\theta' = \tau\tilde{\beta}$. When the classifier is well trained by minimizing the empirical loss, and $\tilde{\delta}_J$ is sufficiently small, the proof can be completed.

S2 Additional numerical analysis results

S2.1 Imbalanced annotations

In practice, it is common for each annotator to label only a subset of the data rather than the entire dataset. For instance, in our real data application on lung disease, each subject receives a limited number of diagnoses from a pool of doctors. Motivated by this setting, we extend our model to address the imbalanced annotator problem.

Suppose there are R annotators, where each annotator r labels a subset of data denoted by $I_r \subset [n]$, with $|I_r| = n_r$ for $r \in [R]$. We adapt the same procedure used in the original minimization problem (3.4) to accommodate this imbalanced setting. The resulting optimization problem is given by

$$(\widehat{\mathbf{f}}^*, \widehat{\mathbf{A}}^{(r)}) = \arg \min_{\mathbf{f} \in \mathcal{D}, \mathbf{A}^{(r)} \in \mathcal{A}} -\frac{1}{R} \sum_{r=1}^R \frac{1}{n_r} \sum_{i \in I_r} \sum_{k=1}^K \mathbb{I}(\tilde{Y}_i^{(r)} = k) \log([\mathbf{A}^{(r)} \mathbf{f}(\boldsymbol{\xi}_J^{(i)})]_k) + \lambda \text{tr} \left(\sum_{r=1}^R \mathbf{A}^{(r)} \right).$$

We adopt the same data-generating strategies used in Models 1 – 8, and randomly assign the index sets $\{I_r\}_{r \in [R]}$ such that $\bigcup_{r \in [R]} I_r = [n]$, and $n_r = \lceil 0.5n \rceil$ for all $r \in [R]$.

This design ensures that each annotator evaluates no more than half of the data. Below, we present Figures S1 to S4, which display the classification accuracy and F-norm error of the confusion matrices for all methods, evaluated under varying sampling frequencies across Models 1 to 8, based on the imbalanced design described above. Consistent with the results from the balanced design, our proposed method, **afDNN**, consistently outperforms all baseline methods. Notably, across all levels of contamination, **afDNN** demonstrates performance comparable to that of the oracle classifier, highlighting its effectiveness in accurately estimating the CM and correctly classifying new subjects. Specifically, the proposed **afDNN** maintains strong performance under low contamination levels ($\epsilon = 0.1, 0.3$), as measured by both classification accuracy and CM estimation. When the noise level increases to $\epsilon = 0.5$, our method still outperforms the alternatives, though it remains below the oracle model. This decline is mainly due to the increasing bias in CM estimation and the sparsity of annotations per annotator.

S2.2 Additional Results for Models 5 to 12

Below, Figures S5 to S11 depict classification accuracy and F-norm error in Models 5 to 12.

S3 ADNI dataset

The PET images underwent spatial normalization and comprehensive post-processing to ensure consistency across subjects. The AD group, which had between three and six clinical visits, contributed scans from their third visit, while the EMCI cohort, with

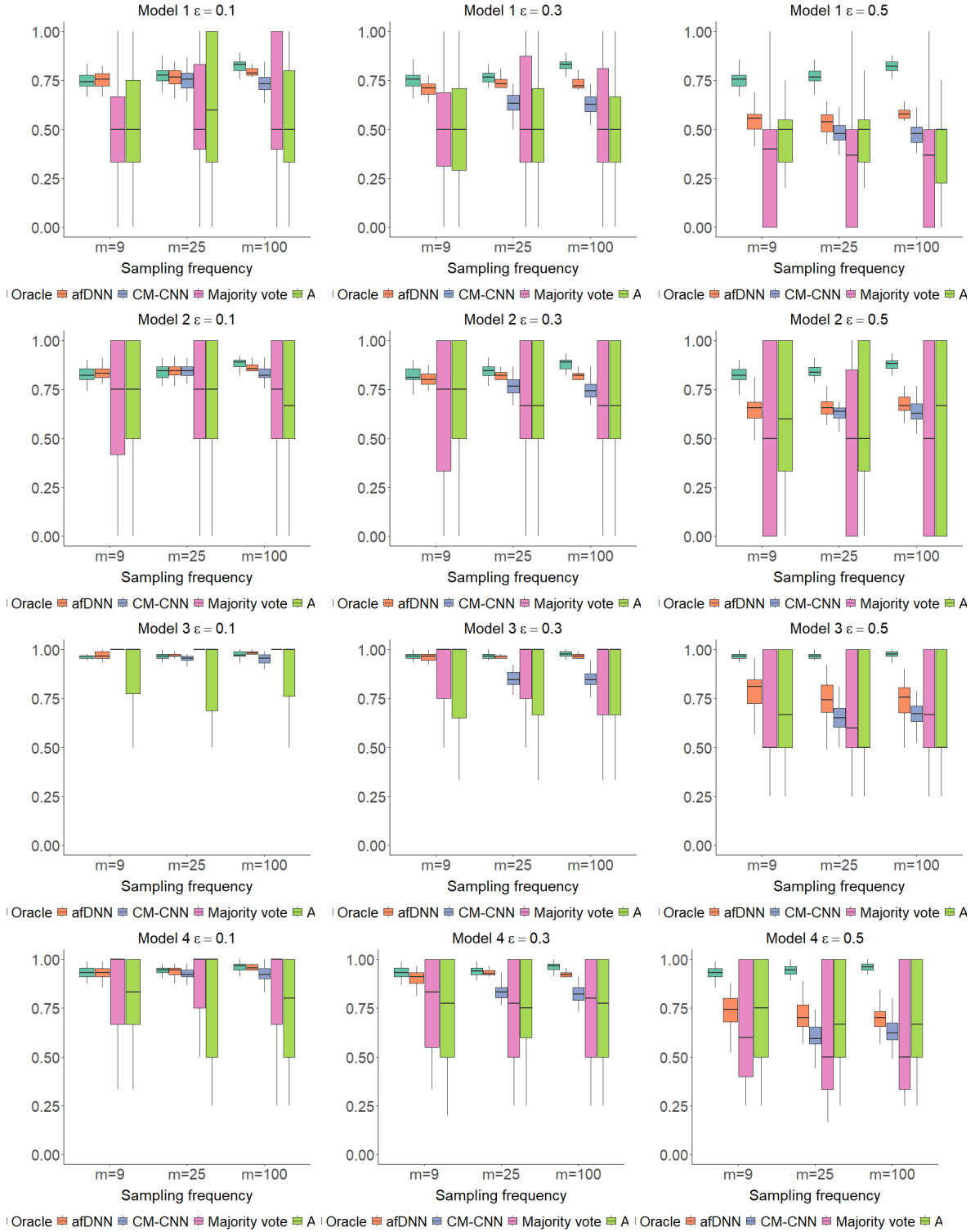


Figure S1: Classification accuracy for all methods across different sampling frequencies for Models 1 – 4 with imbalanced design. From top to bottom, the rows represent the results for Model 1 – 4, respectively. From left to right, the columns correspond to $\epsilon = 0.1, 0.3, \text{ and } 0.5$.

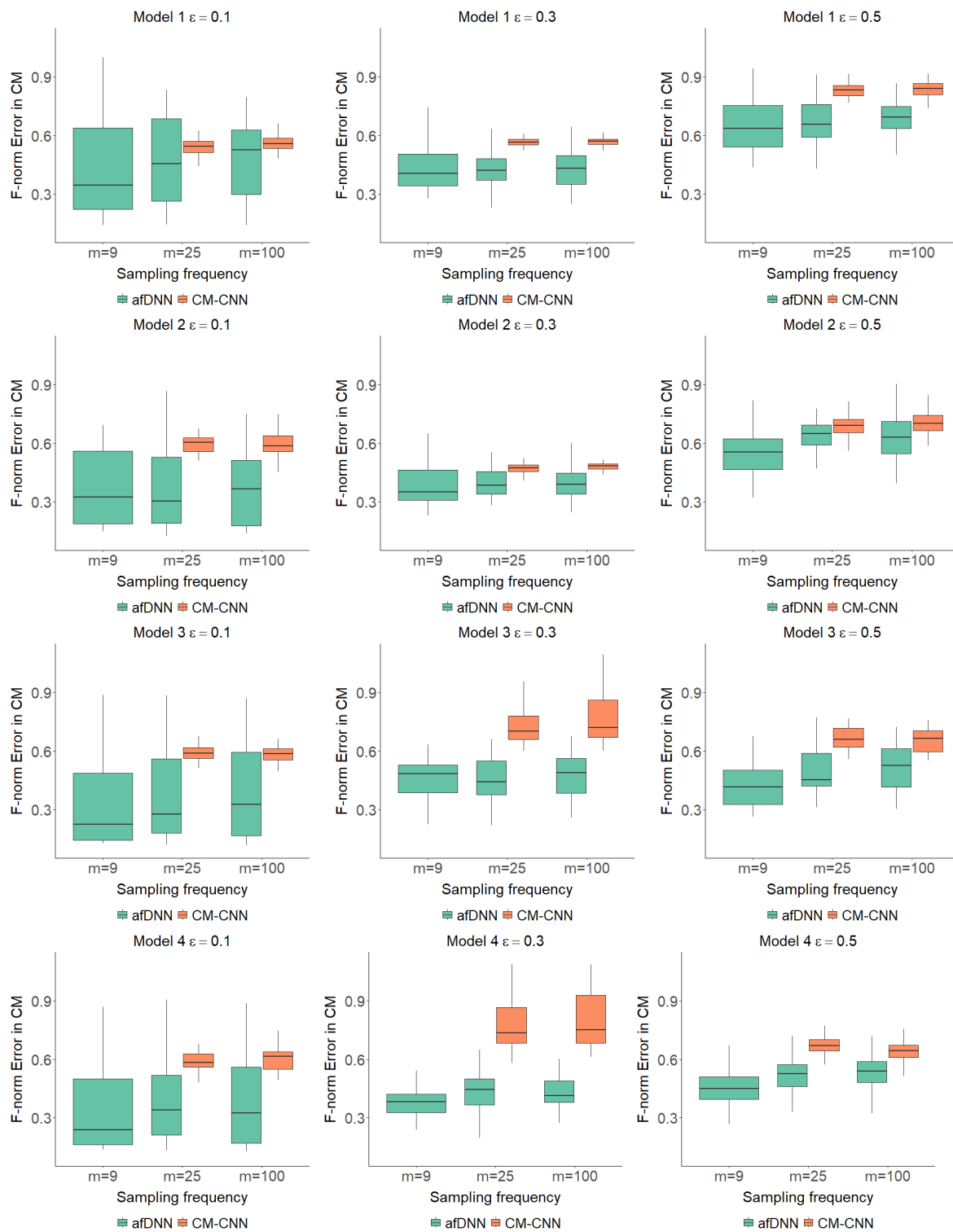


Figure S2: F-norm error in confusion matrices for all methods across different sampling frequencies for Models 1 – 4 with imbalanced design. From top to bottom, the rows represent the results for Model 1 – 4, respectively. From left to right, the columns correspond to $\epsilon = 0.1, 0.3,$ and 0.5 .

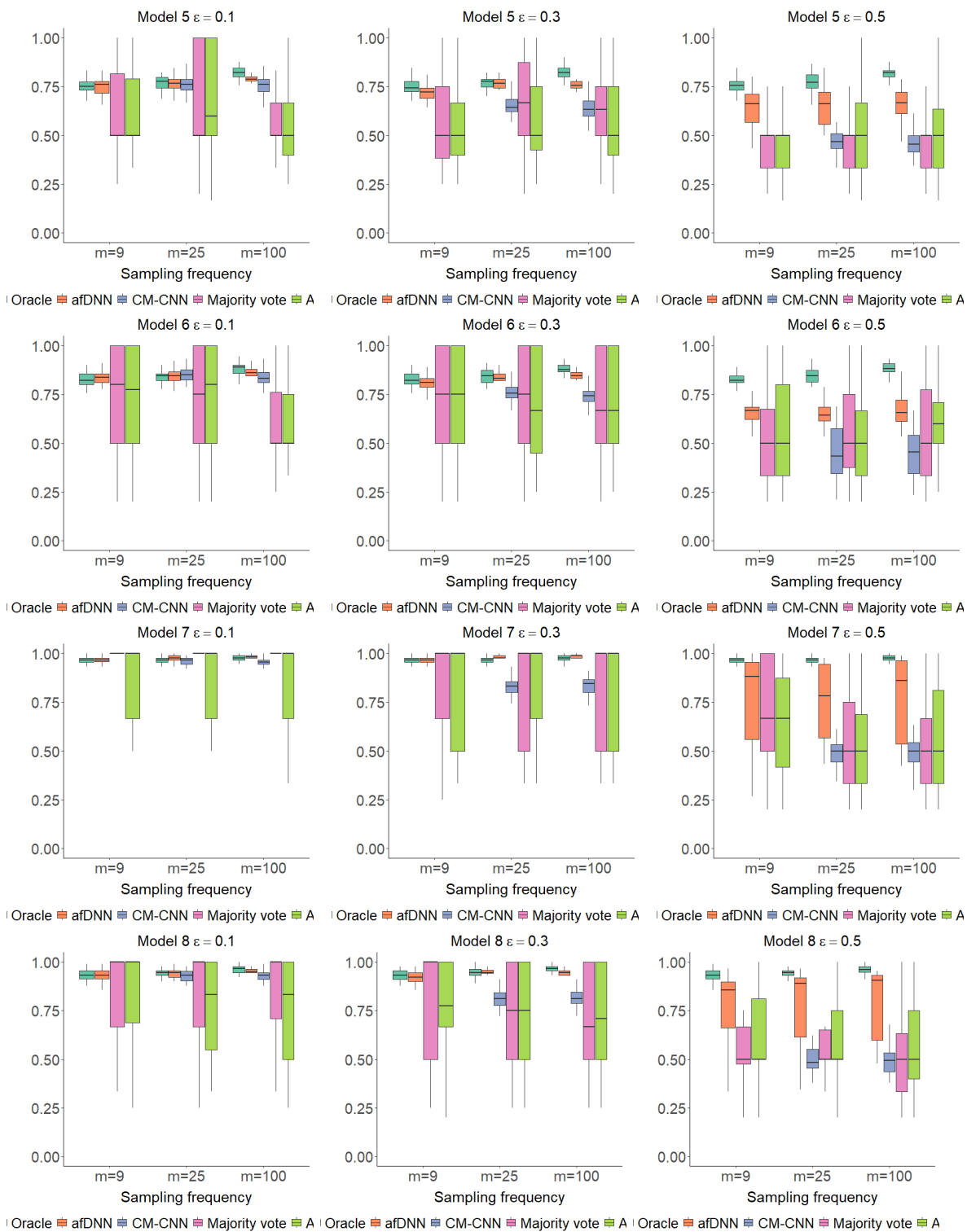


Figure S3: Classification accuracy for all methods across different sampling frequencies for Models 5 – 8 with imbalanced design. From top to bottom, the rows represent the results for Model 5 – 8, respectively. From left to right, the columns correspond to $\epsilon = 0.1, 0.3, \text{ and } 0.5$.

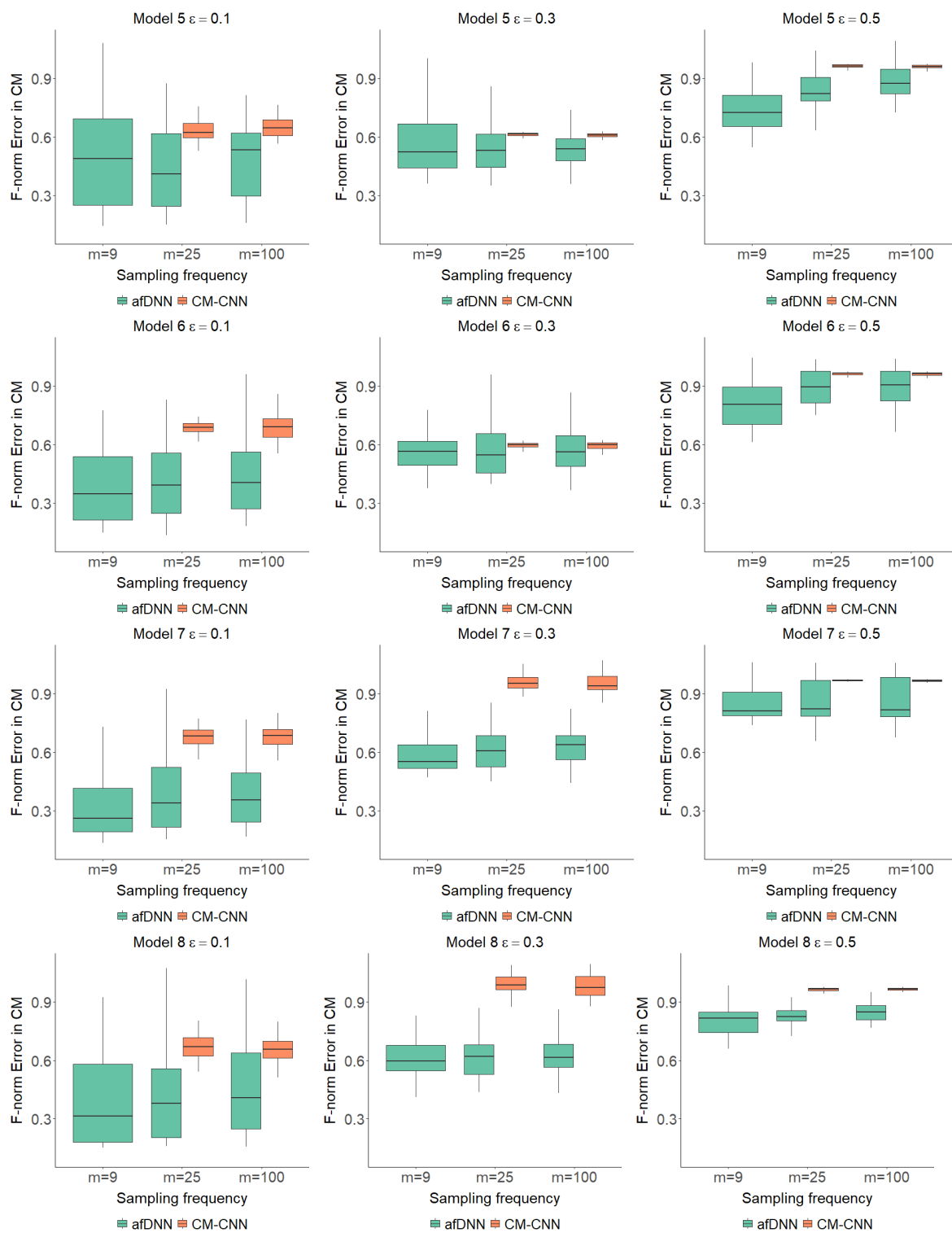


Figure S4: F-norm error in confusion matrices for all methods across different sampling frequencies for Models 5 – 8 with imbalanced design. From top to bottom, the rows represent the results for Model 5 – 8, respectively. From left to right, the columns correspond to $\epsilon = 0.1, 0.3,$ and 0.5 .

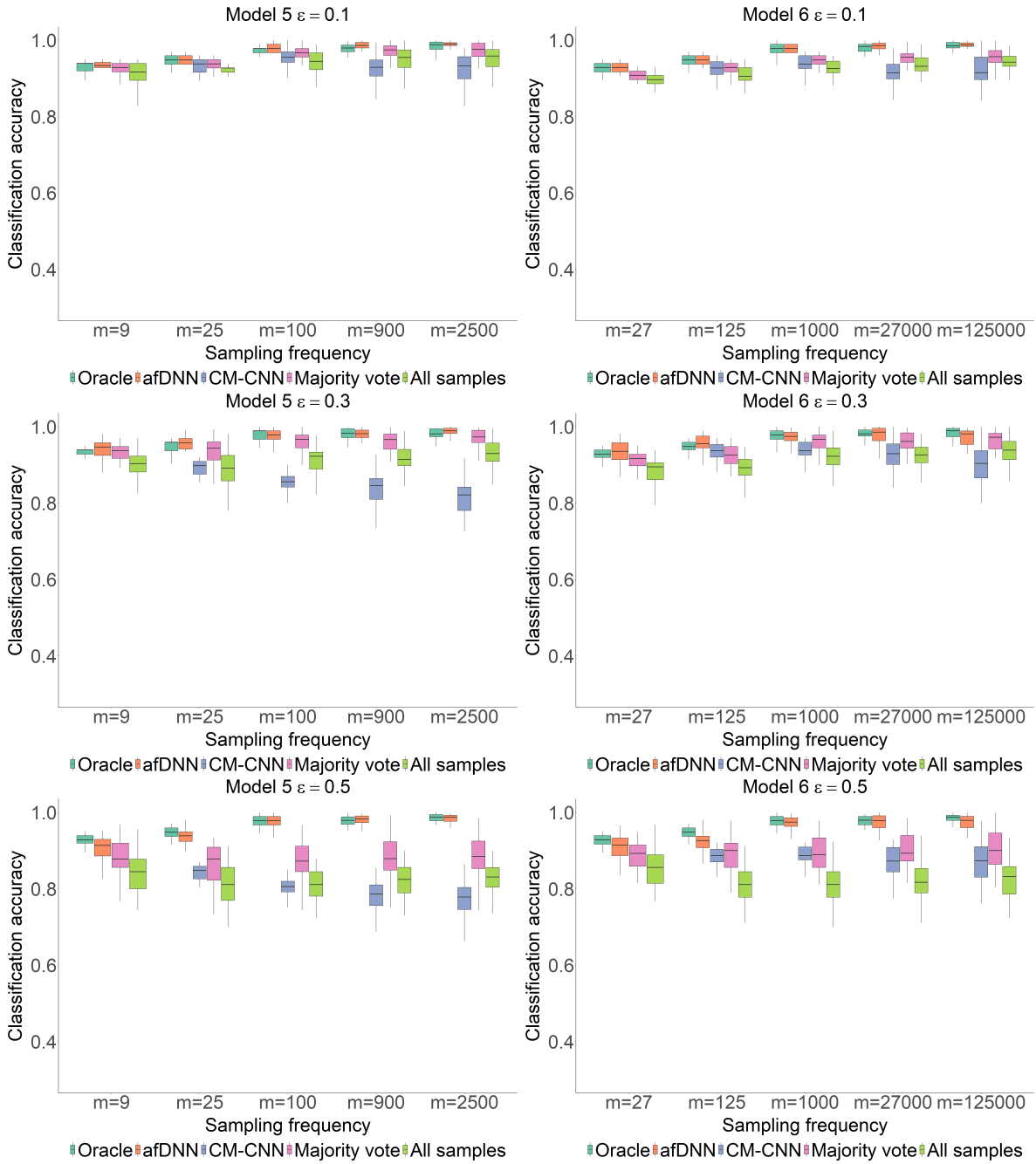


Figure S5: Classification accuracy for all methods across different sampling frequencies for Models 5 – 6. From left to right, the rows represent the results for Model 5 – 6, respectively. From top to bottom, the columns correspond to $\epsilon = 0.1, 0.3$, and 0.5 .

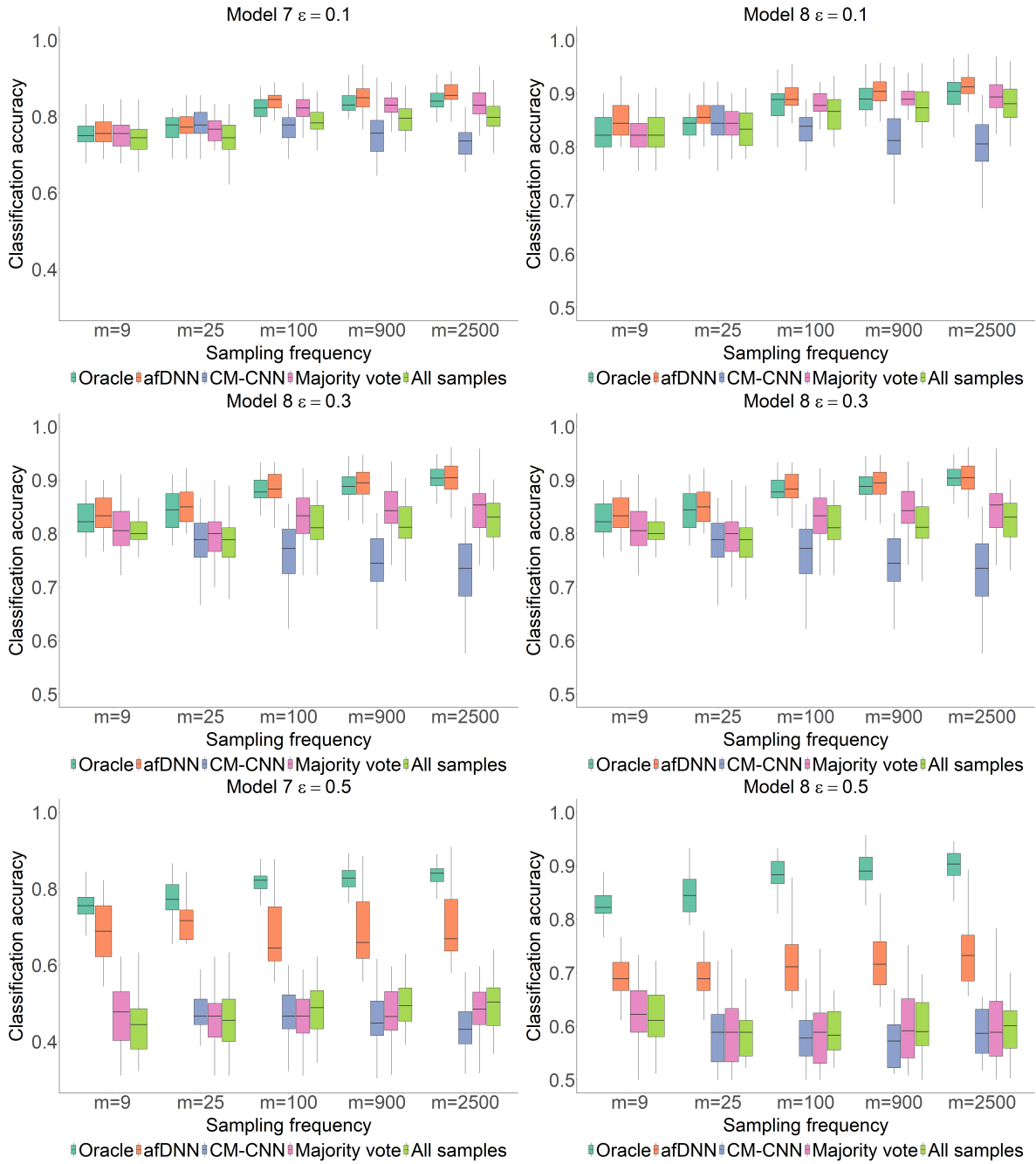


Figure S6: Classification accuracy for all methods across different sampling frequencies for Models 7 – 8. From left to right, the rows represent the results for Model 7 – 8, respectively. From top to bottom, the columns correspond to $\epsilon = 0.1, 0.3$, and 0.5 .

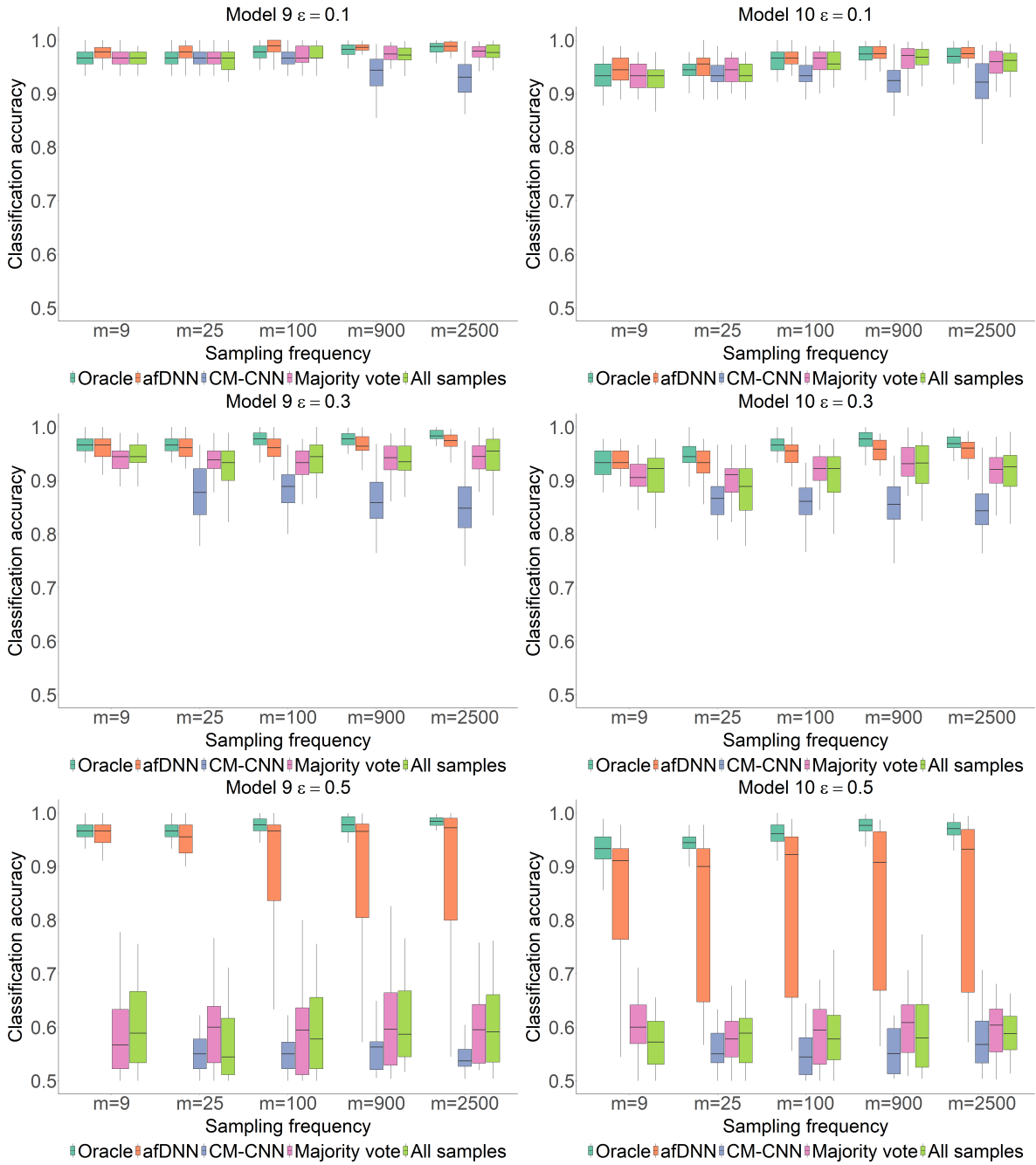


Figure S7: Classification accuracy for all methods across different sampling frequencies for Models 9 – 10. From left to right, the rows represent the results for Model 9 – 10, respectively. From top to bottom, the columns correspond to $\epsilon = 0.1, 0.3$, and 0.5 .

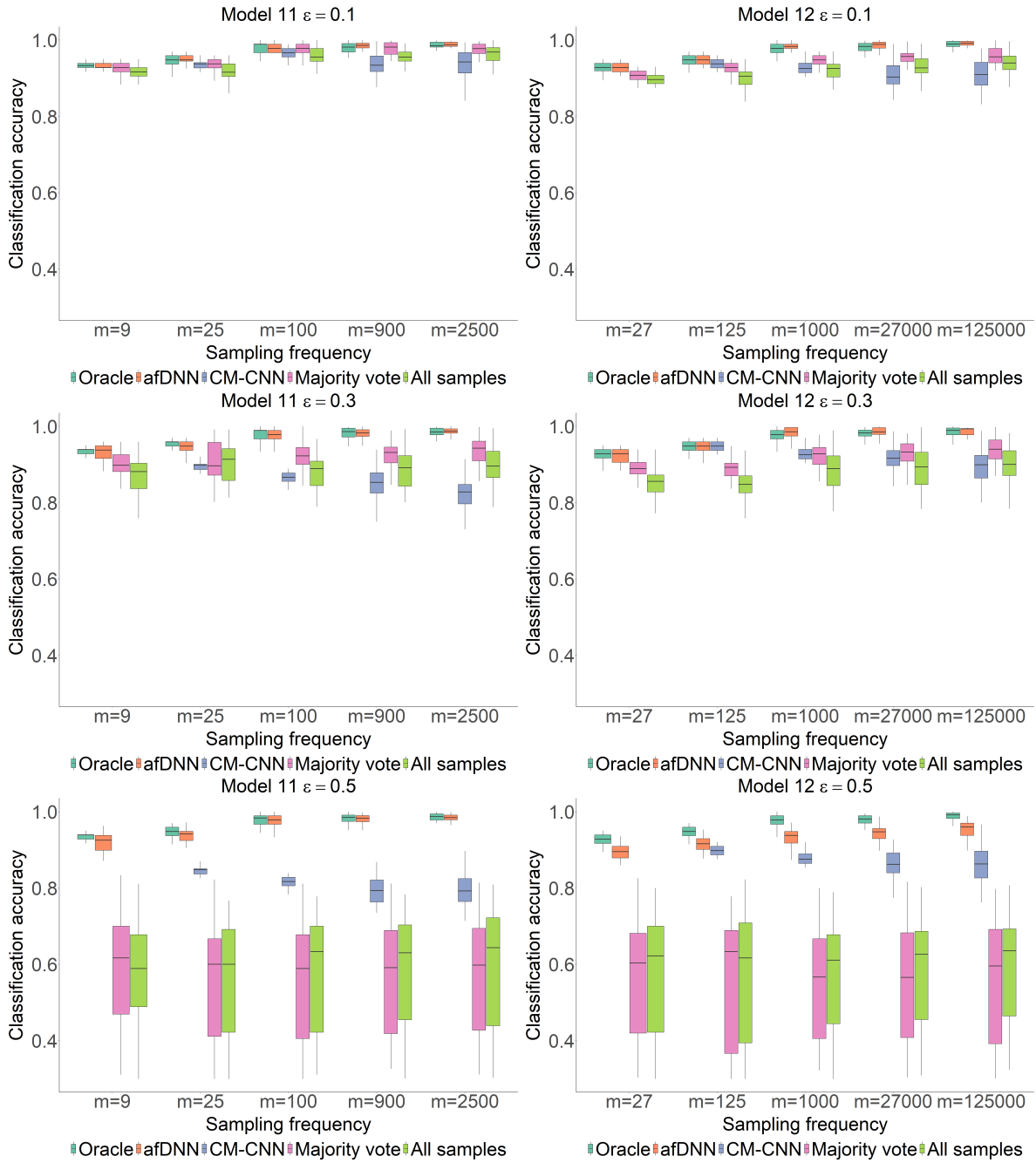


Figure S8: Classification accuracy for all methods across different sampling frequencies for Models 11 – 12. From left to right, the rows represent the results for Model 11 – 12, respectively. From top to bottom, the columns correspond to $\epsilon = 0.1, 0.3$, and 0.5 .

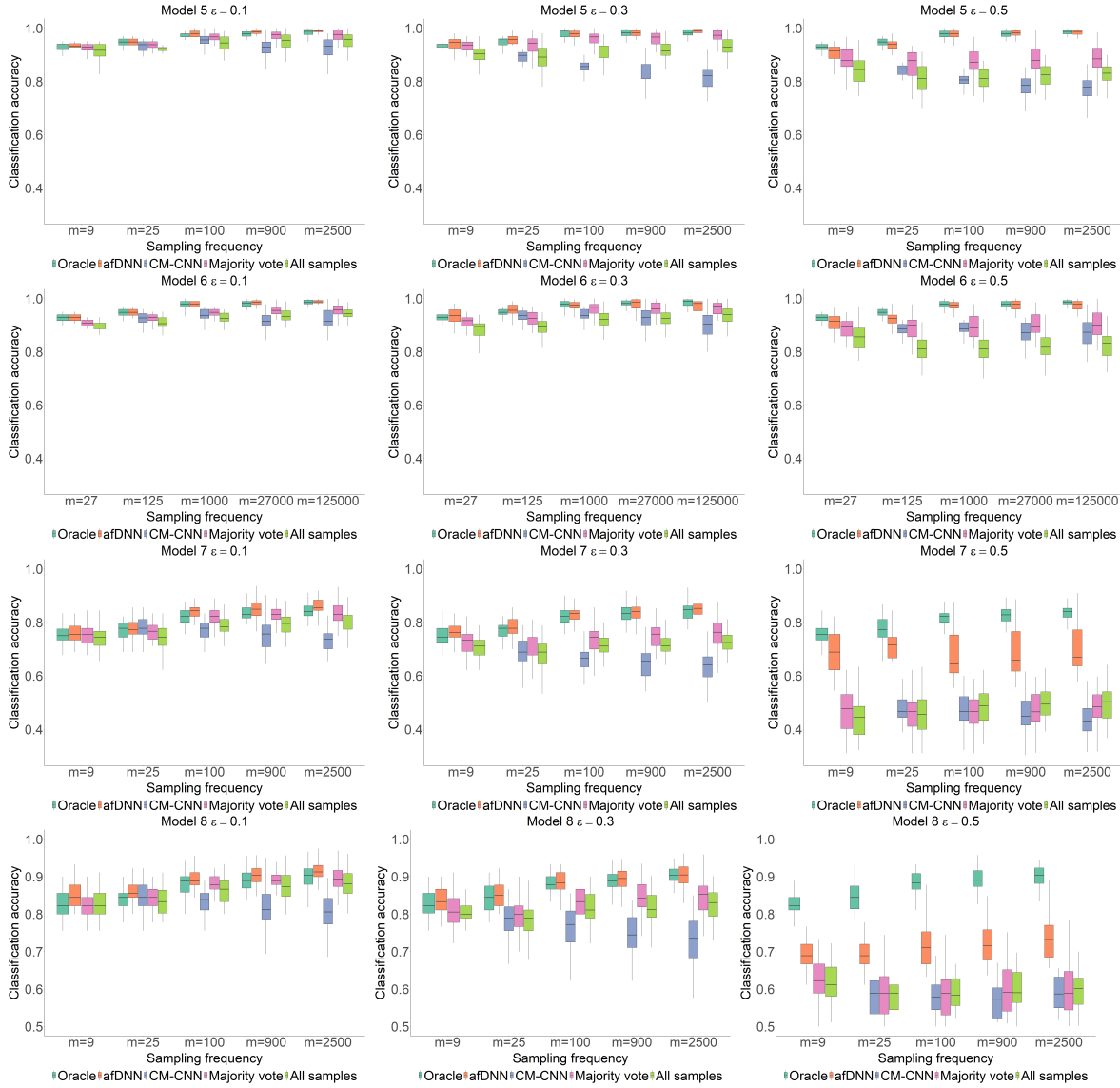


Figure S9: Classification accuracy for all methods across different sampling frequencies for Models 5 – 8. From top to bottom, the rows represent the results for Model 5 – 8, respectively. From left to right, the columns correspond to $\epsilon = 0.1, 0.3$, and 0.5 .

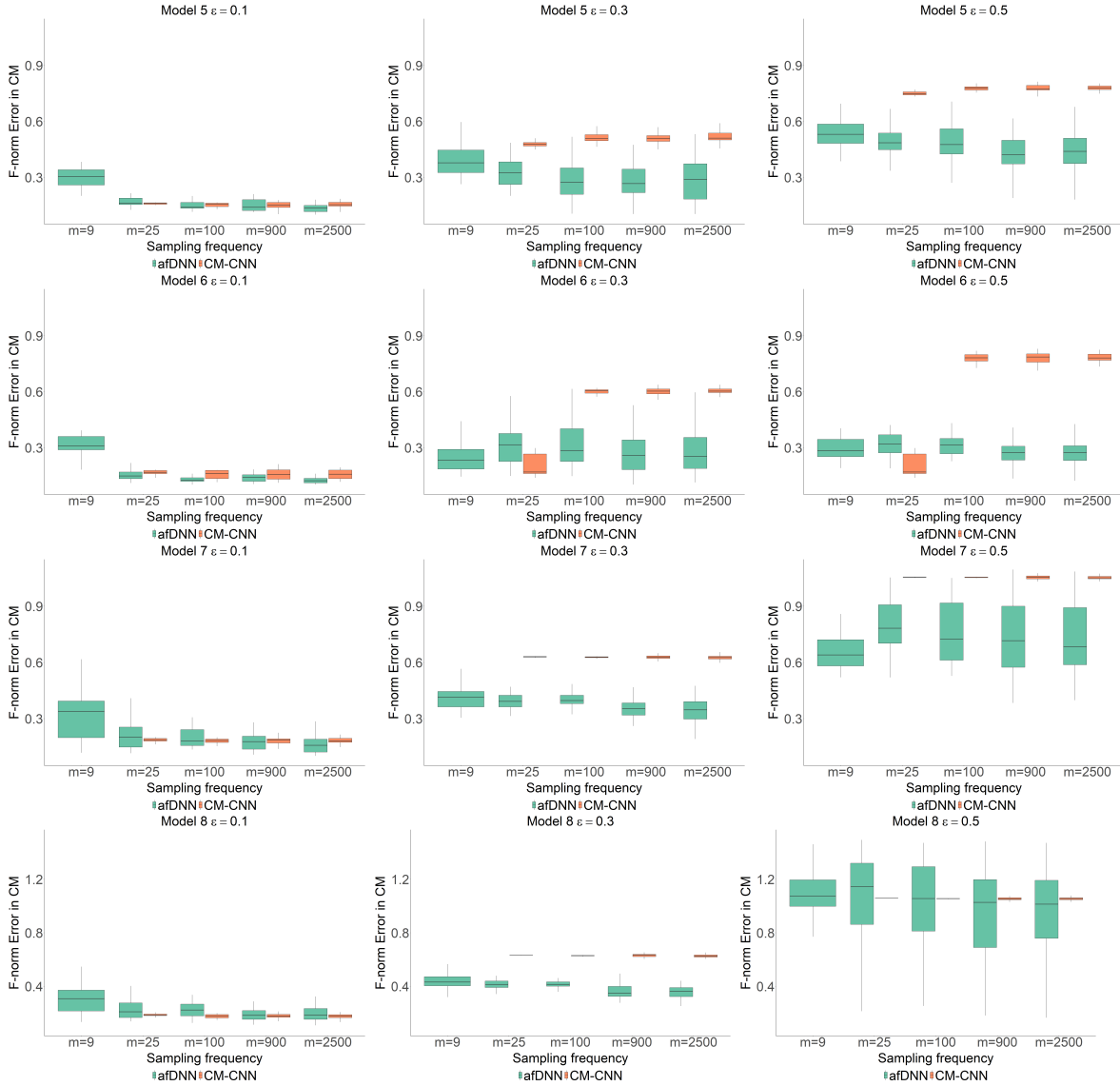


Figure S10: F-norm error in CM for all methods across different sampling frequencies for Models 5 – 8. From top to bottom, the rows represent the results for Model 5 – 8, respectively. From left to right, the columns correspond to $\epsilon = 0.1, 0.3, \text{ and } 0.5$.

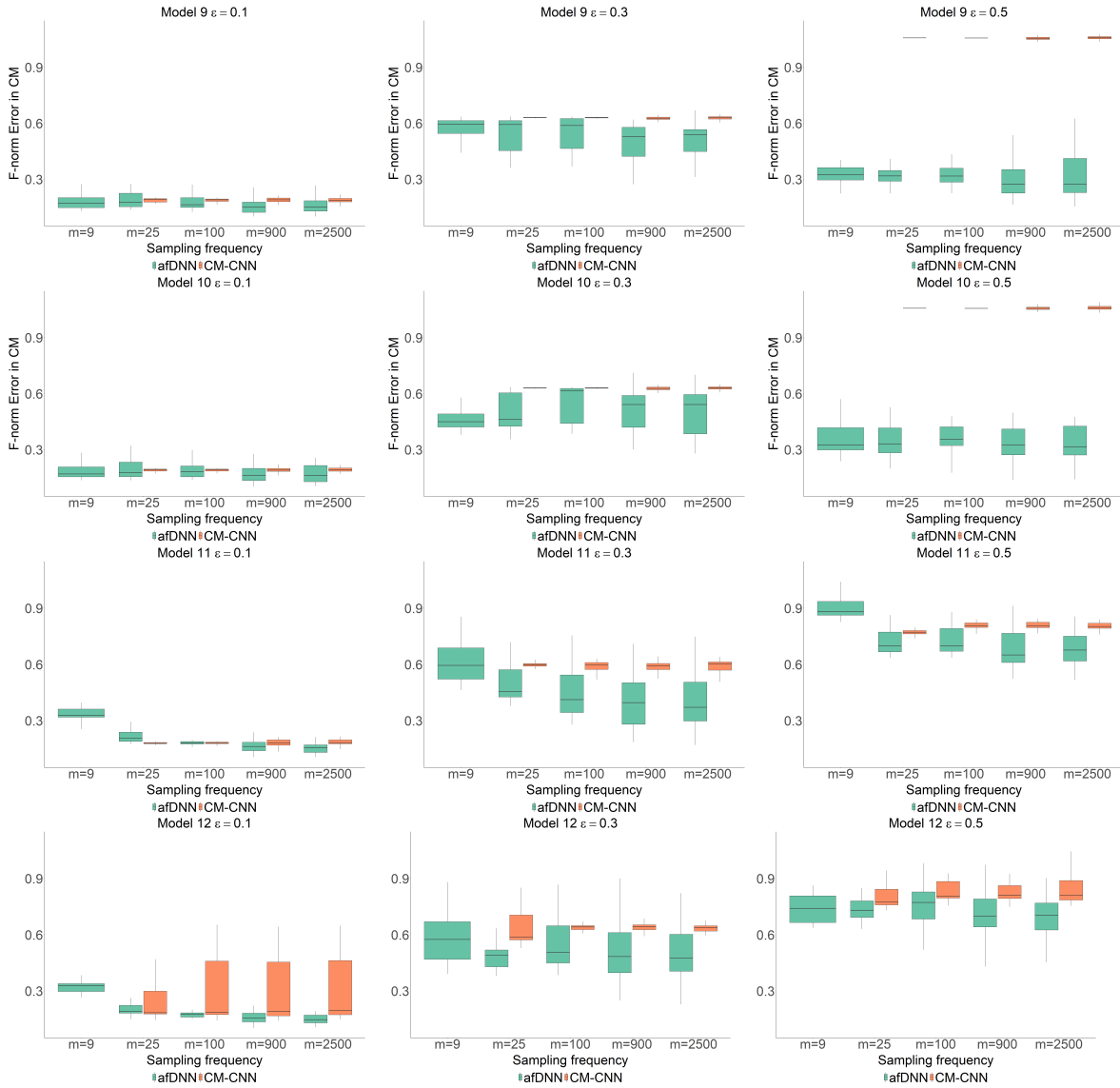


Figure S11: F-norm error in CM for all methods across different sampling frequencies for Models 9 – 12. From top to bottom, the rows represent the results for Model 9 – 12, respectively. From left to right, the columns correspond to $\epsilon = 0.1, 0.3$, and 0.5 .

Table S1: Averaged classification accuracy with standard errors in brackets for ADNI 2D brain images with different values of ϵ for 10th, 30th and 50th slices.

Methods	$\epsilon = 0.1$			$\epsilon = 0.3$			$\epsilon = 0.5$		
	10-th	30-th	50-th	10-th	30-th	50-th	10-th	30-th	50-th
Oracle	0.50 (0.03)	0.45 (0.03)	0.48 (0.04)	0.50 (0.03)	0.45 (0.03)	0.48 (0.04)	0.50 (0.03)	0.45 (0.03)	0.48 (0.04)
afDNN	0.52 (0.02)	0.49 (0.01)	0.52 (0.01)	0.51 (0.02)	0.49 (0.02)	0.51 (0.01)	0.45 (0.03)	0.48 (0.02)	0.46 (0.02)
All sample	0.49 (0.02)	0.47 (0.03)	0.49 (0.01)	0.50 (0.01)	0.48 (0.02)	0.40 (0.01)	0.44 (0.03)	0.43 (0.03)	0.37 (0.03)
Majority vote	0.48 (0.03)	0.45 (0.01)	0.42 (0.02)	0.40 (0.02)	0.47 (0.04)	0.37 (0.03)	0.32 (0.02)	0.36 (0.00)	0.36 (0.00)
CM-CNN	0.50 (0.06)	0.46 (0.06)	0.45 (0.05)	0.43 (0.05)	0.44 (0.06)	0.44 (0.06)	0.38 (0.06)	0.42 (0.06)	0.41 (0.06)

only two visits, provided scans from the second session. Figure S12 illustrates the averaged images of the 20th, 40th and 60th slices for the AD, EMCI, and CN groups. Table S1 presents averaged classification accuracy with standard errors in brackets with different values of ϵ for 10th, 30th and 50th slices. Similar conclusions can be drawn for these three slices, as presented in the main paper: the proposed afDNN consistently demonstrates strong performance under low contamination levels ($\epsilon = 0.1, 0.3$), as reflected in both classification accuracy and confusion matrix estimation. As the noise level increases to $\epsilon = 0.5$, our method continues to outperform competing approaches, although it remains slightly below the oracle model. Table S2 presents the classification accuracy comparison results. All methods achieve slightly better classification accuracy in the 3D setting than the 2D slice-based analysis, suggesting that the additional spatial information is beneficial. Our proposed method continues to perform strongly and remains competitive with the oracle model that is trained using the true labels.

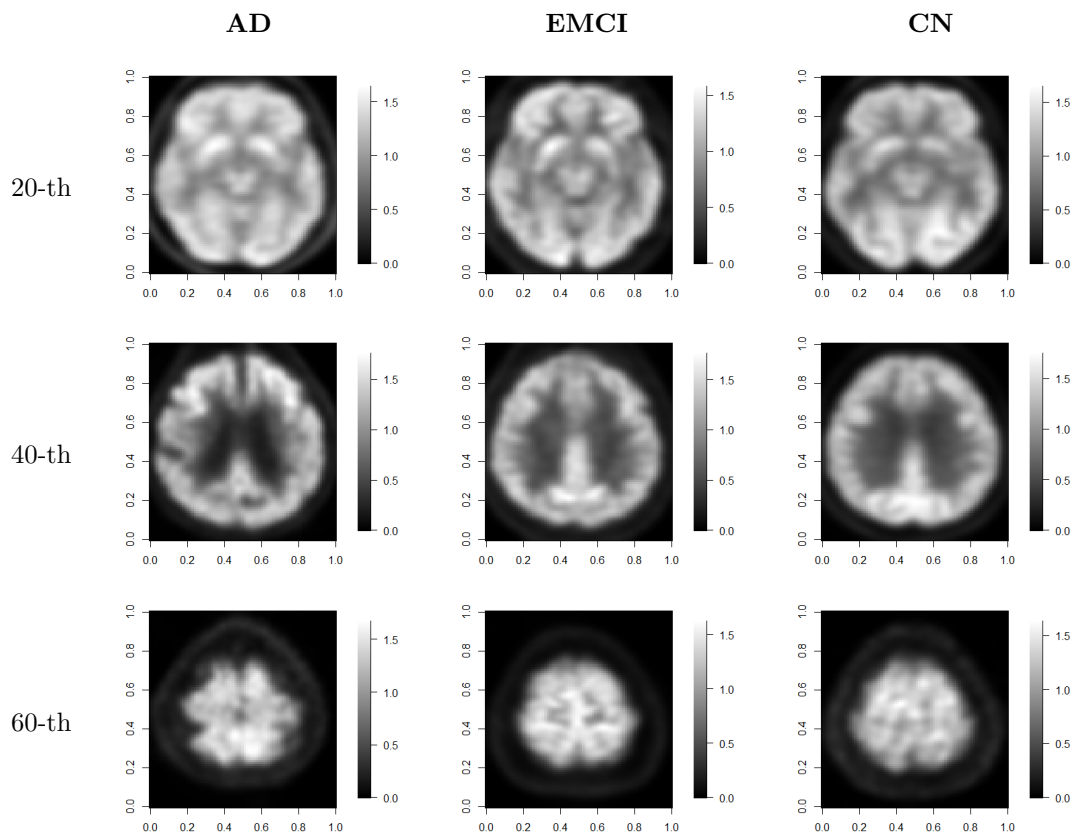


Figure S12: Sample images of the 20th, 40th and 60th slices for AD group (left column), EMCI group (middle column), and CN group (right column).

Table S2: Averaged classification accuracy with standard errors in brackets for ADNI 3D brain images with different values of ϵ .

Methods	$\epsilon = 0.1$	$\epsilon = 0.3$	$\epsilon = 0.5$
Oracle	0.66 (0.07)	0.66 (0.07)	0.66 (0.07)
afDNN	0.66 (0.06)	0.63 (0.06)	0.59 (0.07)
All sample	0.58 (0.05)	0.54 (0.04)	0.50 (0.04)
Majority vote	0.57 (0.05)	0.52 (0.05)	0.51 (0.04)
CM-CNN	0.60 (0.08)	0.58 (0.06)	0.57 (0.06)

Bibliography

- Bos, T. and J. Schmidt-Hieber (2022). Convergence rates of deep relu networks for multiclass classification. *Electronic Journal of Statistics* 16, 2724–2773.
- Tanno, R., A. Saeedi, S. Sankaranarayanan, C. D. Alexander, and N. Silberman (2019). Learning from noisy labels by regularized estimation of annotator confusion. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 11244–11253.
- Wang, S. and G. Cao (2024). Multiclass classification for multidimensional functional data through deep neural networks. *Electronic Journal of Statistics* 18, 1248–1292.