

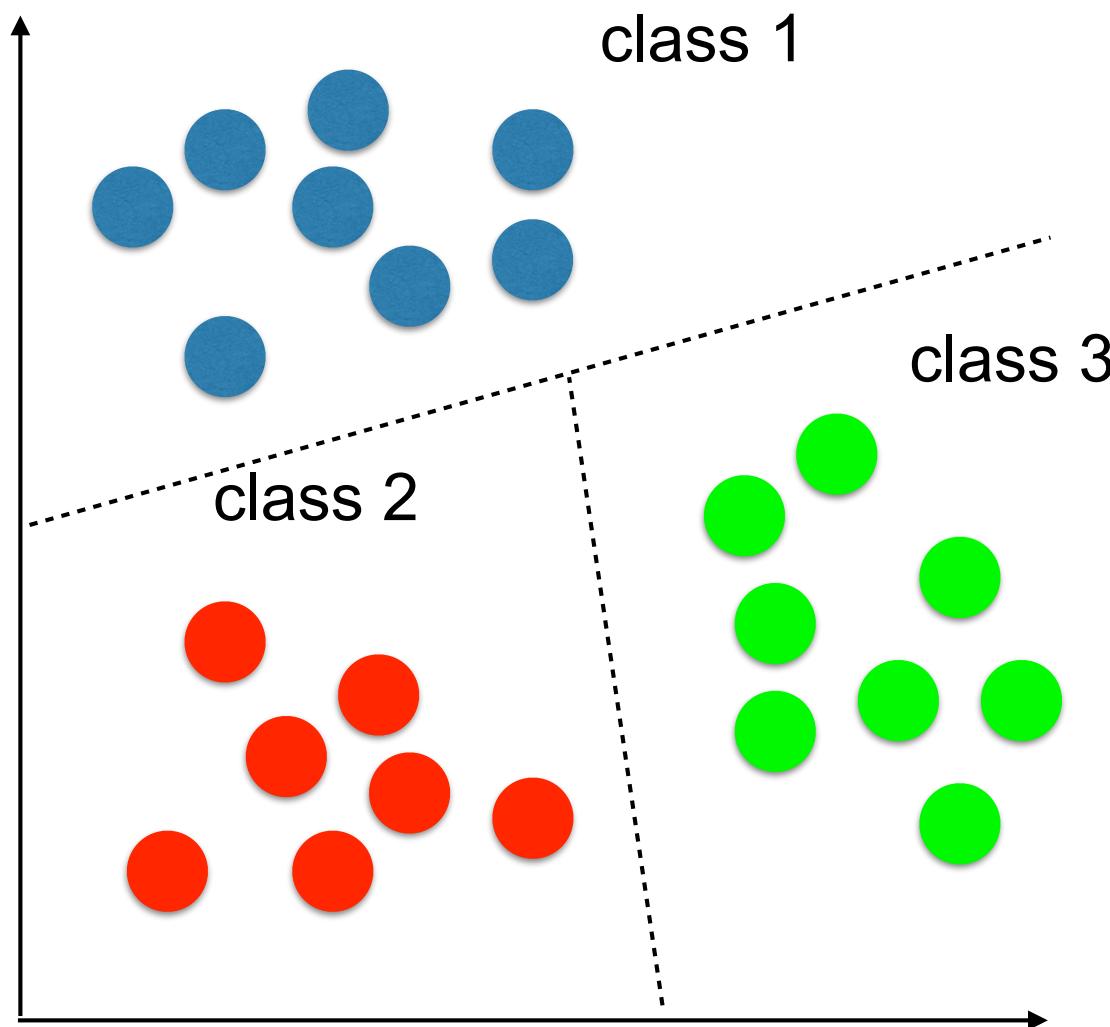
Accelerating Cross-Validation in Multinomial Logistic Regression with ℓ_1 -Regularization

Tomoyuki Obuchi
Dep. Math. Comp. Sci.
Tokyo Tech.

In Collaboration with Yoshiyuki Kabashima in Tokyo Tech.

Multi-class classification

- Under supervised setting



Multinomial Logistic Regression (MLR)

- Multinomial logistic regression (called also as softmax function)

$$P(y = a|x; \{\mathbf{w}_a\}_a) = \frac{e^{\mathbf{x} \cdot \mathbf{w}_a}}{\sum_{b=1}^L e^{\mathbf{x} \cdot \mathbf{w}_b}}$$

- Teacher Data: $D^M = \{(\mathbf{x}_\mu, y_\mu)\}_{\mu=1}^M$
 - Feature vector: $\mathbf{x} \in \mathbb{R}^N$
 - Class variable: $y \in \{1, 2, \dots, L\}$
- A common learning method: Maximum likelihood (ML) approach

$$\{\hat{\mathbf{w}}_a\}_a = \arg \min_{\{\mathbf{w}_a\}_a} \left\{ - \sum_{\mu=1}^M \log \phi \left(y_\mu | \{u_{\mu b} = \mathbf{x}_\mu^\top \mathbf{w}_b\}_b \right) \right\}$$

$\{\mathbf{w}_a \in \mathbb{R}^N\}_{a=1}^L$: Weight or Coefficient

Regularization

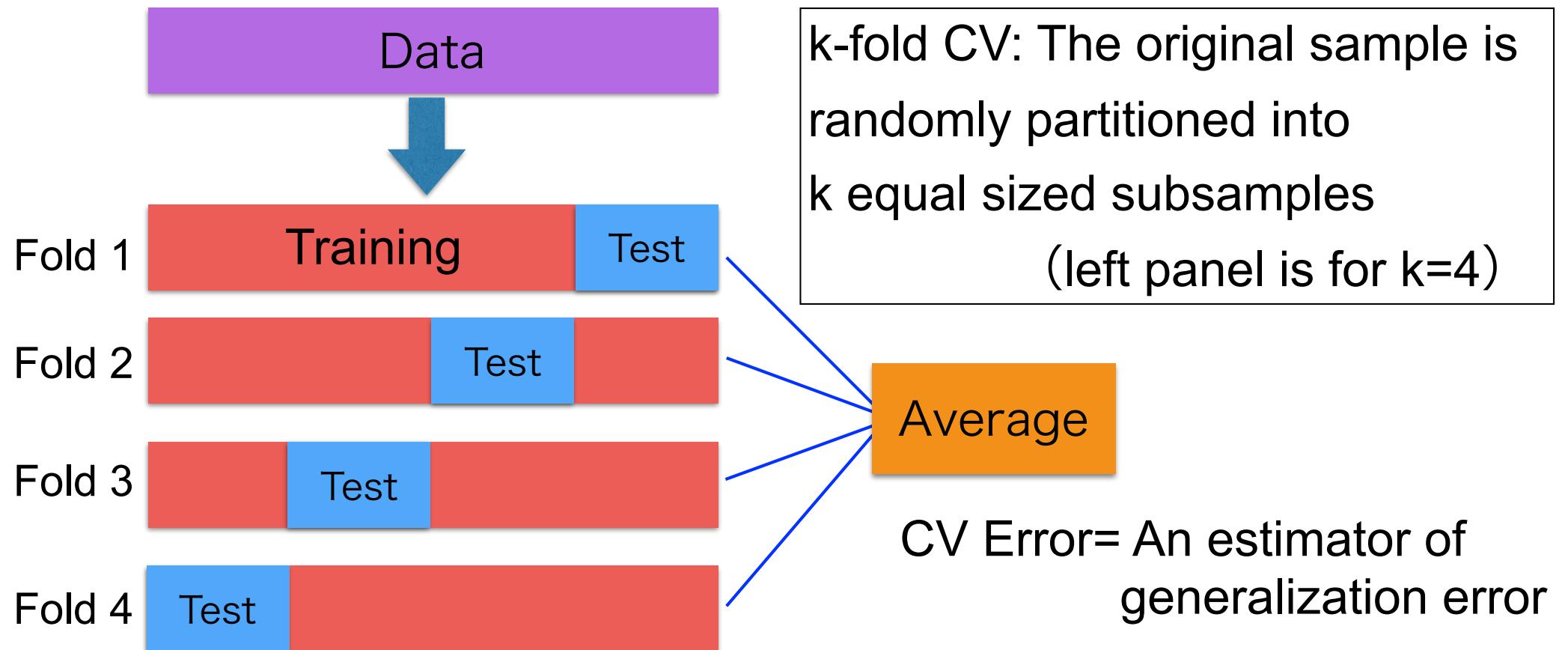
- ML approach can be unstable
 - If the dataset size is not large enough
 - The data generating process largely deviates from the MLR
←Regularization technique
- Penalized ML approach

$$\{\hat{w}_a(\boldsymbol{\lambda})\}_a = \min_{\{\mathbf{w}_a\}_a} \left\{ \sum_{\mu=1}^M q_\mu \left(\{\mathbf{w}_a\}_{a=1}^L \right) + J(\{\mathbf{w}\}_{a=1}^L | \boldsymbol{\lambda}) \right\}$$

- Negative log-likelihood
$$q_\mu \left(\{\mathbf{w}_a\}_{a=1}^L \right) = -\ln \phi \left(y_\mu \mid \left\{ u_{\mu a} = \mathbf{x}_\mu^\top \mathbf{w}_a \right\}_{a=1}^L \right),$$
- How to determine $\boldsymbol{\lambda}$? ←Examine generalization error
→ Cross-Validation (CV)

Cross-validation (CV)

- Dividing dataset into training set and test set
 - Learning on Training Set
 - Examine generalization error on Test Set



LOOCV and Linear Response

- Leave-one-out (LOO) CV (=M-fold CV), Predictive likelihood

$$\{\hat{\mathbf{w}}_a^{\setminus \mu}(\lambda)\}_a = \arg \min_{\{\mathbf{w}_a\}_a} \left\{ \mathcal{H} \left(\{\mathbf{w}_a\}_{a=1}^L \middle| D^M, \lambda \right) - q_\mu \left(\{\mathbf{w}_a\}_{a=1}^L \right) \right\}$$

$$\epsilon_{\text{LOO}}(\lambda) = \frac{1}{M} \sum_{\mu=1}^M q_\mu \left(\left\{ \hat{\mathbf{w}}_a^{\setminus \mu} \right\}_{a=1}^L \right) \quad \text{← Computationally expensive (M times repeated optimization)}$$

- Diff. between Full and LOO solutions owes to only one term of q_μ → Linear (response) approximation!
- Approximate formula (called ACV formula):

$$\hat{\mathbf{u}}_\mu^{\setminus \mu} \approx \hat{\mathbf{u}}_\mu + C_\mu (I - C_\mu F_\mu)^{-1} \mathbf{b}_\mu$$

- Crucial assumption for the formula: The set of active variables is common between full and LOO systems.

Ingredients in approximation

$$\hat{\mathbf{u}}_\mu^{\setminus \mu} \approx \hat{\mathbf{u}}_\mu + C_\mu (I - C_\mu F_\mu)^{-1} \mathbf{b}_\mu$$

$$C_\mu = X_{*\hat{A}}^\mu (G_{\hat{A}\hat{A}})^{-1} \left(X_{*\hat{A}}^\mu \right)^\top \quad G_{(ai)(bj)} \equiv \frac{\partial^2 \mathcal{H}}{\partial w_{ai} \partial w_{bj}} \equiv (\partial^2 \mathcal{H})_{(ai)(bj)}$$

$$\hat{A} = \{(ai) | \hat{w}_{ai} \neq 0\}$$

$$X^\mu = \underbrace{\begin{pmatrix} \boxed{x_\mu^\top} & \cdot & \cdot & \cdot & 0 \\ 0 & \ddots & & & \boxed{x_\mu^\top} \end{pmatrix}}_{NL} \Biggr\}^L$$

$$F_{ab}^\mu \equiv \frac{\partial^2 q_\mu}{\partial u_{\mu a} \partial u_{\mu b}} = \delta_{ab} p_{a|\mu} - p_{a|\mu} p_{b|\mu} \quad p_{a|\mu} = \phi(a | \{\hat{u}_{\mu b}\}_b)$$

$$\mathbf{b}_\mu = \frac{\partial q_\mu}{\partial \mathbf{u}_\mu} = (p_{1|\mu} - \delta_{1y_\mu}, p_{2|\mu} - \delta_{2y_\mu}, \dots, p_{L|\mu} - \delta_{Ly_\mu})^\top$$

Derivation of ACV

- Equation determining the weights

$$\nabla \mathcal{H}(\hat{\mathbf{w}}) = 0 \Rightarrow \hat{\mathbf{w}}$$

$$\nabla \mathcal{H}^{\backslash \mu}(\hat{\mathbf{w}}) = 0 \Rightarrow \hat{\mathbf{w}}^{\backslash \mu}$$

- Perturbation assuming the smallness of $d_\mu = \hat{\mathbf{w}} - \hat{\mathbf{w}}^{\backslash \mu}$

$$0 = \nabla \mathcal{H}^{\backslash \mu}(\hat{\mathbf{w}}^{\backslash \mu}) = \nabla \mathcal{H}^{\backslash \mu}(\hat{\mathbf{w}} - d_\mu) \approx \nabla \mathcal{H}^{\backslash \mu}(\hat{\mathbf{w}}) - \partial^2 \mathcal{H}^{\backslash \mu}(\hat{\mathbf{w}}) d_\mu$$

$$= \nabla(\mathcal{H}(\hat{\mathbf{w}}) - q_\mu(\hat{\mathbf{w}})) - \partial^2 \mathcal{H}^{\backslash \mu}(\hat{\mathbf{w}}) d_\mu$$

$$= -\nabla q_\mu(\hat{\mathbf{w}}) - G^{\backslash \mu}(\hat{\mathbf{w}}) d_\mu,$$

$$\Rightarrow d_\mu \approx -\left(G^{\backslash \mu}(\hat{\mathbf{w}})\right)^{-1} \nabla q_\mu(\hat{\mathbf{w}}).$$

$$G^{\backslash \mu}(\hat{\mathbf{w}}) = G(\hat{\mathbf{w}}) - \partial^2 q_\mu(\hat{\mathbf{w}}),$$

$$\Rightarrow \boxed{\hat{\mathbf{w}}^{\backslash \mu} = \hat{\mathbf{w}} - d_\mu \approx \hat{\mathbf{w}} + (G(\hat{\mathbf{w}}) - \partial^2 q_\mu(\hat{\mathbf{w}}))^{-1} \nabla q_\mu(\hat{\mathbf{w}})}$$

Derivation of ACV

- Multiplying x_μ^\top and simplifying factors by Sherman-Morrison formula

$$\hat{w}^{\setminus \mu} = \hat{w} - d_\mu \approx \hat{w} + (G(\hat{w}) - \partial^2 q_\mu(\hat{w}))^{-1} \nabla q_\mu(\hat{w})$$

$$\Rightarrow \hat{u}_\mu^{\setminus \mu} \approx \hat{u}_\mu + C_\mu (I - C_\mu F_\mu)^{-1} b_\mu$$

- To apply the derivation so far to the case with sparse regularizations, the computation of gradient and Hessian is done only on the active set $\hat{A} = \{(ai) | \hat{w}_{ai} \neq 0\}$
 - We also need to assume that the active set is common between full and LOO systems
 - This is ``approximately'' true
 - Negligible contributions by variables come and go support

Warning

- The factor $C_\mu = X_{*\hat{A}}^\mu (G_{\hat{A}\hat{A}})^{-1} \left(X_{*\hat{A}}^\mu \right)^\top$ is *not well defined*
 - Because the Hessian G has zero eigenvalues
 - The model is singular because the MLR is invariant against
$$\mathbf{w}_a \rightarrow \mathbf{w}_a + \mathbf{v} \ (\forall a)$$
 - This is gauge degree of freedom (gdf) in physics
 - Some implementation erases this gdf by fixing certain class weights at zeros (Krishnapuram et al., 2005; Schmidt, 2010)
 - This is nice for our approximation
 - But not all implementations do this
 - E.g. *Glmnet* (Friedman et al., 2010) doesn't

→ Another way should be tailored to avoid the problem

Zero Mode Removal by Hand

- The factor $C_\mu = X_{*\hat{A}}^\mu (G_{\hat{A}\hat{A}})^{-1} \left(X_{*\hat{A}}^\mu \right)^\top$ is *not well defined*
 - Because the Hessian G has zero eigenvalues
 - Remove the zero modes by hand

$$G_{\hat{A}\hat{A}} \stackrel{\text{EVD}}{=} \sum_i d_i \mathbf{v}_i \mathbf{v}_i^\top = \sum_{i \in S^+} d_i \mathbf{v}_i \mathbf{v}_i^\top,$$


$$\overline{G}_{\hat{A}\hat{A}}^{-1} \equiv \sum_{i \in S^+} d_i^{-1} \mathbf{v}_i \mathbf{v}_i^\top.$$


$$C_\mu = X_{*\hat{A}}^\mu \overline{G}_{\hat{A}\hat{A}}^{-1} \left(X_{*\hat{A}}^\mu \right)^\top.$$

This seems to be ad-hoc but actually works!

Algorithm and Cost (ACV)

Algorithm 1 Approximate CV of the MLR

```
1: procedure ACV( $\hat{W}(\lambda), D^M$ )
2:   Compute the active set  $\hat{A}$  from  $\hat{W}$ 
3:   Compute  $\{\hat{u}_\mu, X^\mu, b^\mu, F^\mu\}_\mu$  by eqs. (1,7), (12) and (13)
4:    $G_{\hat{A}\hat{A}} \leftarrow \sum_{\mu=1}^M (X^\mu)^\top F^\mu X^\mu$                                  $\triangleright O(ML|\hat{A}|^2 + ML^2|\hat{A}|)$ 
5:   Compute  $\bar{G}_{\hat{A}\hat{A}}^{-1}$  by eq. (25)                                               $\triangleright O(|\hat{A}|^3)$ 
6:   for  $\mu = 1, \dots, M$  do
7:      $C_\mu \leftarrow X_{*\hat{A}}^\mu \bar{G}_{\hat{A}\hat{A}}^{-1} (X_{*\hat{A}}^\mu)^\top$ 
8:      $\hat{u}_\mu^\backslash \leftarrow \hat{u}_\mu + C_\mu (I_L - F^\mu C_\mu)^{-1} b^\mu$ 
9:   end for
10:  Compute  $\epsilon_{\text{LOO}}$  from  $\{u_\mu^\backslash\}_\mu$  by eq. (6)
11:  return  $\epsilon_{\text{LOO}}$ 
12: end procedure
```

- Comp. cost = $O(ML|\hat{A}|^2 + |\hat{A}|^3) \approx O(MN^2 + N^3)$
 - Too heavy if N, M are large

Self-Averaging Approximation (SAACV)

- Neglecting ``correlations'' bet. different feature components

$$\left(G^{\setminus \mu}\right)^{-1}_{(ai)(bj)} \approx \begin{cases} (\chi_i)_{ab} \delta_{ij}, & ((ai), (bj) \in \hat{A}) \\ 0, & (\text{otherwise}) \end{cases},$$

$\chi_i \in \mathbb{R}^{L \times L}$: Inter-class susceptibility (rescaled variance)

- Self-consistent equation to determine χ

$$\begin{cases} (\chi_i)_{\hat{A}_i \hat{A}_i} = \frac{1}{\sigma_x^2} \left(\sum_{\nu=1}^M \left((I_L + F^\nu C_{\text{SA}})^{-1} F^\nu \right)_{\hat{A}_i \hat{A}_i} \right)^{-1} \\ C_{\text{SA}} = \sigma_x^2 \sum_{i=1}^N \chi_i \end{cases}$$

$\hat{A}_i = \{a | \hat{w}_{ai} \neq 0\}$ $\sigma_x^2 = \frac{1}{NM} \sum_{\mu,i} x_{\mu i}^2$

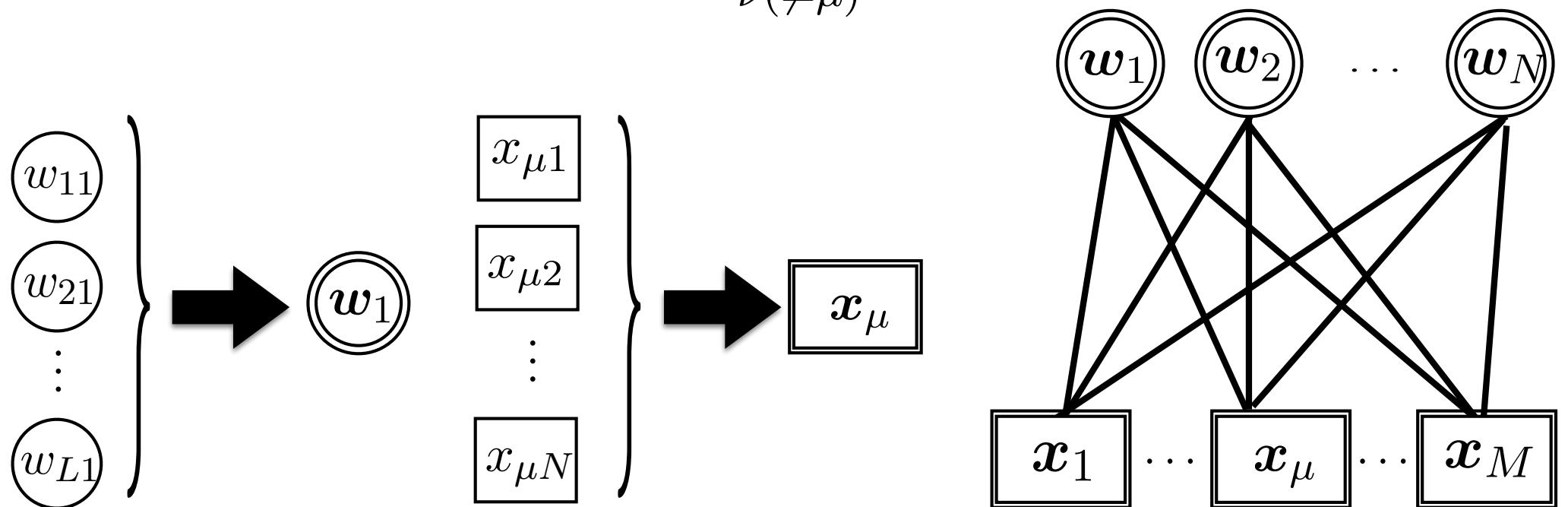
- Computational cost: $O(NL^3 + ML^3)$ Linear in N and M !

Derivation of SAACV

- Cavity method/Belief propagation for multicomponent model

$$\tilde{M}_{\mu \rightarrow i}(\mathbf{w}_i) = \int \prod_{j(\neq i)} d\mathbf{w}_j \phi^\beta(\mathbf{u}_\mu) \prod_{j(\neq i)} M_{j \rightarrow \mu}(\mathbf{w}_j),$$

$$M_{i \rightarrow \mu}(\mathbf{w}_i) = e^{-\beta \lambda ||\mathbf{w}_i||_1} \prod_{\nu(\neq \mu)} \tilde{M}_{\nu \rightarrow i}(\mathbf{w}_i),$$



Derivation of SAACV

- Cavity method/Belief propagation for multicomponent model

$$\tilde{M}_{\mu \rightarrow i}(\mathbf{w}_i) = \int \prod_{j(\neq i)} d\mathbf{w}_j \phi^\beta(\mathbf{u}_\mu) \prod_{j(\neq i)} M_{j \rightarrow \mu}(\mathbf{w}_j),$$

$$u_{\mu a} = \sum_i x_{\mu i} w_{ai} \approx x_{\mu i} w_{ai} + \sum_{j(\neq i)} x_{\mu j} \langle w_{aj} \rangle^\mu + t_a, \quad \mathbf{t} \sim \mathcal{N} \left(0, \text{Cov} \left(\sum_{j(\neq i)} x_{\mu j} \mathbf{w}_j \right) \right)$$

$$\chi_{(ai)(bj)}^{\mu} \equiv \beta \left(\langle w_{ai} w_{bj} \rangle^\mu - \langle w_{ai} \rangle^\mu \langle w_{bj} \rangle^\mu \right) \approx \delta_{ij} \beta \left(\langle w_{ai} w_{bi} \rangle^\mu - \langle w_{ai} \rangle^\mu \langle w_{bi} \rangle^\mu \right)$$

$$\left(C_\mu^{\mu} \right)_{ab} \equiv \sum_{i,j} x_{\mu i} x_{\mu j} \chi_{(ai)(bj)}^{\mu} \approx \sum_i x_{\mu i}^2 \left(\chi_i^{\mu} \right)_{ab}$$

$$\tilde{M}_{\mu \rightarrow i}(\mathbf{w}_i) \approx \int d\mathbf{t} e^{\beta \left(-\frac{1}{2} \mathbf{t}^\top (C_\mu^{\mu})^{-1} \mathbf{t} - q_\mu(\mathbf{w}_i, \mathbf{t}) \right)} \equiv \int d\mathbf{t} e^{\beta f^\mu(\mathbf{w}_i, \mathbf{t})}$$

- The remaining derivation is the same as usual BP

Algorithm and Cost (SAACV)

Algorithm 2 Self-averaging approximate CV of the MLR

```

1: procedure SAACV( $\hat{\mathbf{W}}(\lambda), D^M$ )
2:   Compute the active sets  $\{\hat{A}_i\}_{i=1}^N$  from  $\hat{\mathbf{W}}$ 
3:   Compute  $\{\mathbf{u}_\mu, \mathbf{X}^\mu, \mathbf{b}_\mu, \mathbf{F}^\mu\}_\mu$  by eqs. (1,7), (12) and (13)
4:    $t \leftarrow 0$                                      ▷ Start initialization
5:   for  $i = 1, \dots, N$  do
6:      $(\chi_i^{\setminus\mu})^{(t)} \leftarrow 0,$ 
7:      $(\chi_i^{\setminus\mu})_{\hat{A}_i \hat{A}_i}^{(t)} \leftarrow \sigma_x^{-2},$ 
8:   end for
9:    $\Delta \leftarrow 100$                                 ▷ End initialization
10:  while  $\Delta > \theta$  do                         ▷ Compute  $C_{\text{SA}}$  by recursion
11:     $C_{\text{SA}}^{(t+1)} \leftarrow \sigma_x^2 \sum_{i=1}^N (\chi_i^{\setminus\mu})^{(t)}$ 
12:     $R \leftarrow \sum_{\mu=1}^M (I_L + \mathbf{F}^\mu C_{\text{SA}}^{(t+1)})^{-1} \mathbf{F}^\mu$           ▷  $O(ML^3)$ 
13:     $\Delta \leftarrow 0$ 
14:    for  $i = 1, \dots, N$  do
15:      Compute  $\bar{R}_{\hat{A}_i \hat{A}_i}^{-1}$  by eq. (41) from  $R$ 
16:       $(\chi_i^{\setminus\mu})_{\hat{A}_i \hat{A}_i}^{(t+1)} \leftarrow \sigma_x^{-2} \bar{R}_{\hat{A}_i \hat{A}_i}^{-1}$           ▷  $O(NL^3)$ 
17:       $\Delta \leftarrow \Delta + \|(\chi_i^{\setminus\mu})_{\hat{A}_i \hat{A}_i}^{(t+1)} - (\chi_i^{\setminus\mu})_{\hat{A}_i \hat{A}_i}^{(t)}\|_F$ 
18:    end for
19:     $\Delta \leftarrow \Delta/N$ 
20:     $t \leftarrow t + 1$ 
21:  end while
22:  for  $\mu = 1, \dots, M$  do
23:     $\mathbf{u}_\mu^{\setminus\mu} \leftarrow \mathbf{u}_\mu + C_{\text{SA}}^{(t)} \mathbf{b}^\mu$ 
24:  end for
25:  Compute  $\epsilon_{\text{LOO}}$  from  $\{\mathbf{u}_\mu^{\setminus\mu}\}_\mu$  by eq. (6)
26:  return  $\epsilon_{\text{LOO}}$ 
27: end procedure

```

Simulated Data

- True feature vector is sparse

$$\mathbf{w}_{0a} \sim \prod_{i=1}^N \{(1 - \rho_0)\delta(w_{0ai}) + \rho_0 \mathcal{N}(0, 1/\rho_0)\}$$

- Class label y_μ is drawn uniformly and randomly

$$y_\mu \sim (1/L) \sum_{a=1}^L \delta_{ay_\mu}$$

- Observed Feature=True Feature+Noise

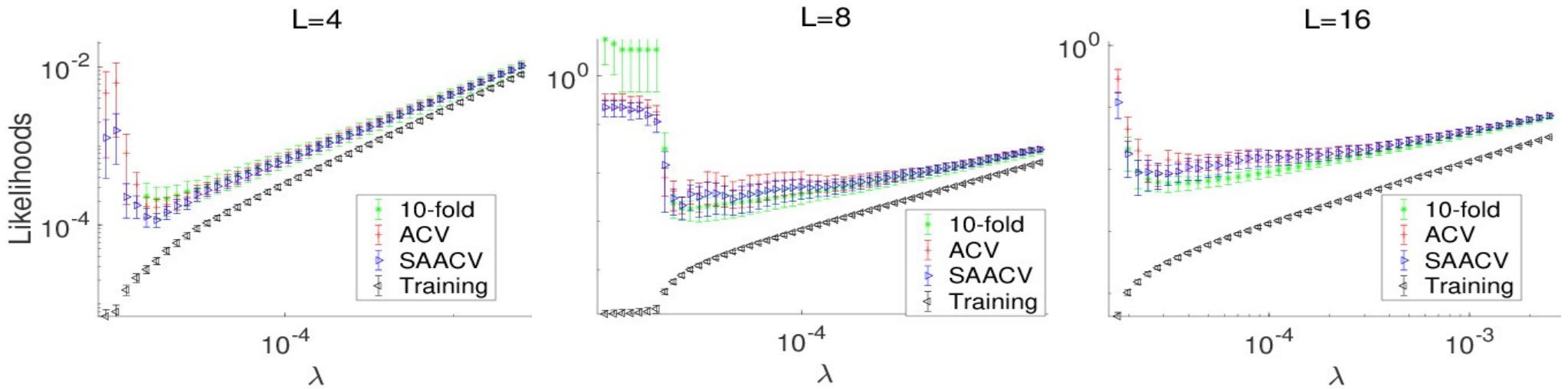
$$x_\mu = \frac{\mathbf{w}_{0y_\mu}}{\sqrt{N}} + \boldsymbol{\xi}, \quad \boldsymbol{\xi}_i \sim \mathcal{N}(0, \sigma_N^2)$$

- The overall parameters= $\{N, L, \alpha, \rho_0, \sigma_\xi^2\}$, $\left(\alpha = \frac{M}{N}\right)$

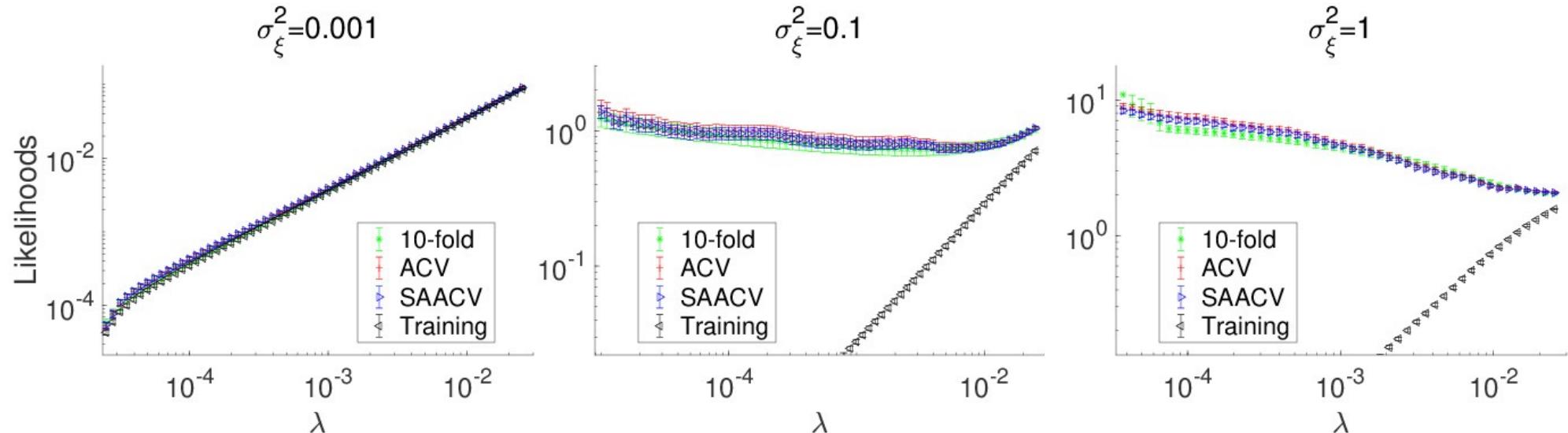
Fix $\alpha = 2, \rho_0 = 0.5$ for simplicity below

Dependence on Class Number and Noise Strength

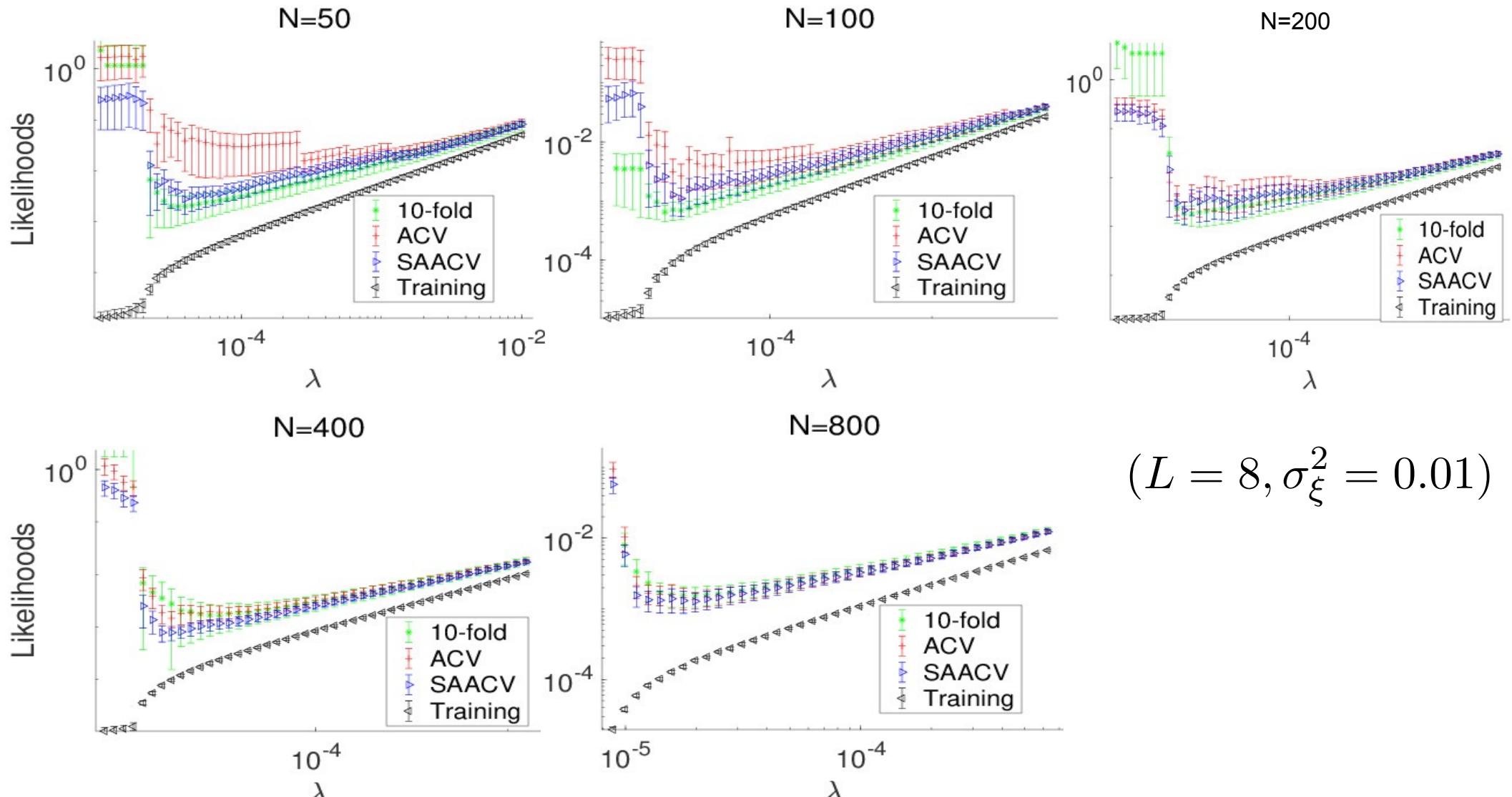
$$(N = 200, \sigma_\xi^2 = 0.01)$$



$$(N = 200, L = 8)$$

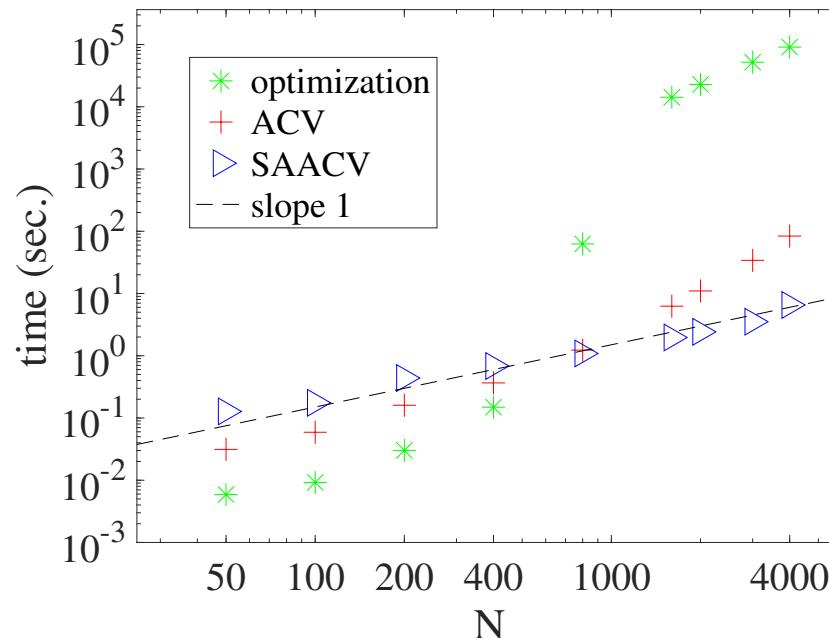


System-Size Dependence



Accuracy seems to be good for $N \geq 200$ even at small λ s.

Computational Time



- Optimization is done by Glmnet
 - A specialized version of coordinate descent method
 - The comp. cost is supposed to be $O(N^2)$
- Small ($N < 400$): Opt < ACV < SAACV \rightarrow Literal CV
- Middle ($400 \leq N < 800$): ACV < SAACV << Opt \rightarrow ACV
- Large ($800 \leq N$): SAACV << ACV << Opt \rightarrow SAACV

ISOLET DATA

ISOLET Data Set

Download: [Data Folder](#), [Data Set Description](#)

Abstract: Goal: Predict which letter-name was spoken—a simple classification task.

Data Set Characteristics:	Multivariate	Number of Instances:	7797	Area:	Computer
Attribute Characteristics:	Real	Number of Attributes:	617	Date Donated	1994-09-12
Associated Tasks:	Classification	Missing Values?	No	Number of Web Hits:	82016

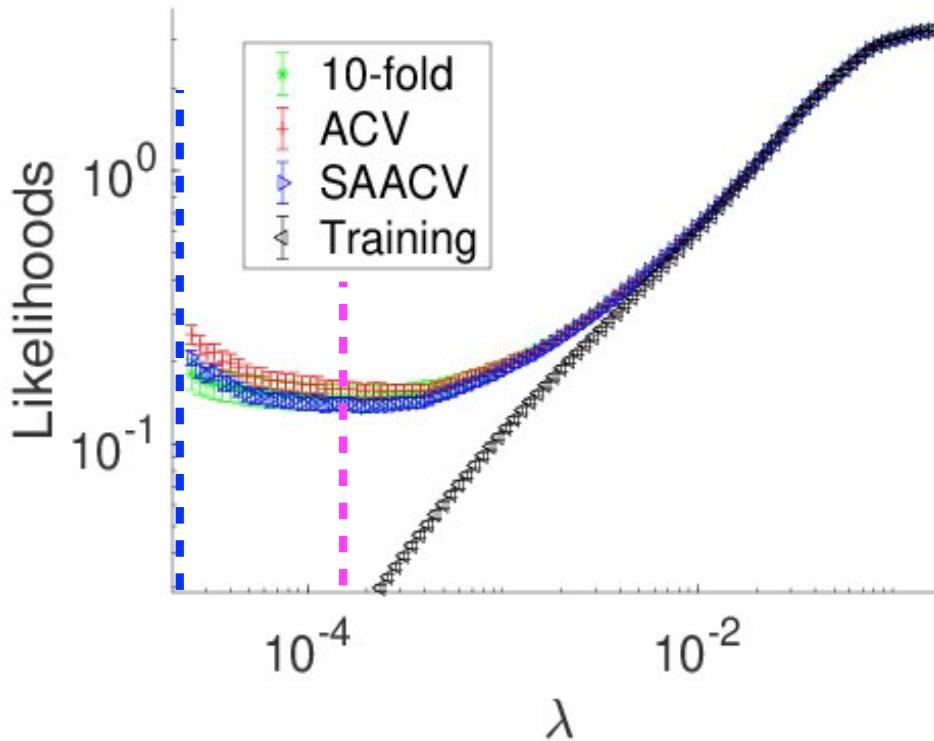
Data Set Information:

This data set was generated as follows. 150 subjects spoke the name of each letter of the alphabet twice. Hence, we have 52 training examples from each speaker. The speakers are grouped into sets of 30 speakers each, and are referred to as isolet1, isolet2, isolet3, isolet4, and isolet5. The data appears in

- Classification task of N=617, M=6238, and L=26

ISOLET DATA

ISOLET data, 26 classes



Actual computational time:

10-fold	ACV	SAACV
7825 s	5173 s	689 s

Accuracy rate at the pred. lik. minimum:

10-fold	ACV	SAACV
0.86	0.86	0.86

Accuracy rate at the minimum λ :

10-fold	ACV	SAACV
0.83	0.78	0.81

- ISOLET: ℓ_1 regularization is meaningful
 - N,M,L are large → ACV is not so good, SAACV is great :)

Summary

[TO, YK: arXiv: 1711.05420](#)

- An approximate formula of CV for the MLR with the ℓ_1 regularization is developed.
 - Perturbation using the largeness of model&dataset
 - Two versions: ACV, SAACV
 - ACV comp. cost: $O(ML|\hat{A}|^2 + |\hat{A}|^3) \approx O(MLN^2 + N^3)$
 - SAACV comp. cost: $O(ML^3 + NL^3)$ (linear in $M\&N$)
 - Cavity method for multicomponent model is used
 - Block-wise correlations are taken into account
 - Numerical experiments demonstrate when these approximations are good and advantageous
 - Matlab, python implementations are available from
 - https://github.com/T-Obuchi/AcceleratedCVonMLR_matlab
 - https://github.com/T-Obuchi/AcceleratedCVonMLR_python
- Elastic net ($\ell_1 + \ell_2$) regularization is also covered in the implementation