

# A Responsible Softmax Layer in Deep Learning

Ryan Coatney

University of Arizona

18 June 2020

# OVERVIEW

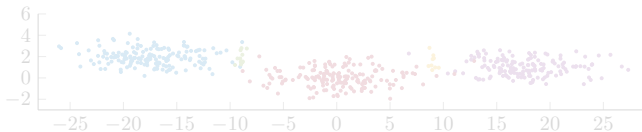
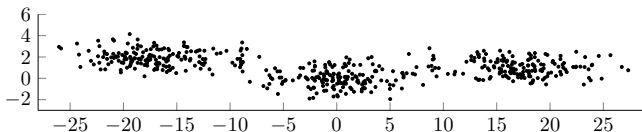
Clustering and Classification

Dynamic Responsibility

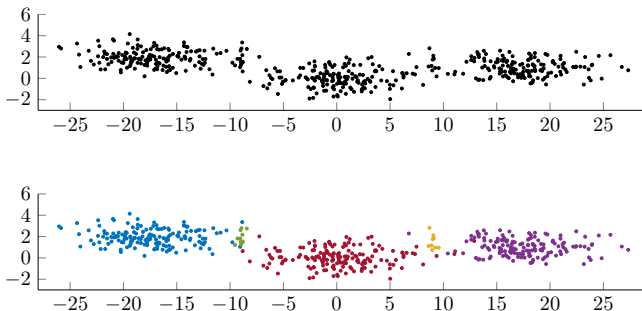
Responsible Softmax

Basic Experiments

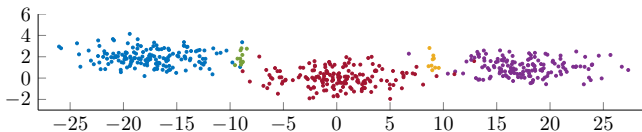
# CLUSTERING



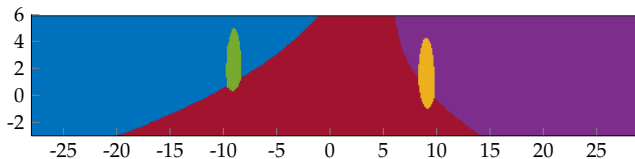
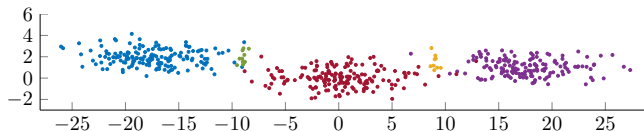
# CLUSTERING



# CLASSIFICATION



# CLASSIFICATION



# K-MEANS ALGORITHM

Start with data  $\mathcal{X} = \{\mathbf{x}^1, \dots, \mathbf{x}^N\}$  and  $K$  starting ‘means’,  $\{\mathbf{m}_1, \dots, \mathbf{m}_K\}$ .

**Assignment:** For each data point,  $\mathbf{x}^n$ , set  $\hat{k}_n = \arg \min_k d(\mathbf{x}^n, \mathbf{m}_k)$ . Set  $\rho_i^n = \delta_i^{\hat{k}_n}$  (Hard responsibility).

**Update:** Let  $N_k = \sum_{n=1}^N \rho_k^n$  and

$$\mathbf{m}_k^{new} = \frac{\sum_{n=1}^N \rho_k^n \mathbf{x}^{(n)}}{N_k}.$$

If  $d(\mathbf{m}_k, \mathbf{m}_k^{new})$  is small for every  $k$ , stop. Otherwise set  $\mathbf{m}_k = \mathbf{m}_k^{new}$  and repeat. See MacKay (2002) for more discussion.

# $K$ -MEANS ALGORITHM

Start with data  $\mathcal{X} = \{\mathbf{x}^1, \dots, \mathbf{x}^N\}$  and  $K$  starting ‘means’,  $\{\mathbf{m}_1, \dots, \mathbf{m}_K\}$ .

**Assignment:** For each data point,  $\mathbf{x}^n$ , set  $\hat{k}_n = \arg \min_k d(\mathbf{x}^n, \mathbf{m}_k)$ . Set  $\rho_i^n = \delta_i^{\hat{k}_n}$  (Hard responsibility).

**Update:** Let  $N_k = \sum_{n=1}^N \rho_k^n$  and

$$\mathbf{m}_k^{new} = \frac{\sum_{n=1}^N \rho_k^n \mathbf{x}^{(n)}}{N_k}.$$

If  $d(\mathbf{m}_k, \mathbf{m}_k^{new})$  is small for every  $k$ , stop. Otherwise set  $\mathbf{m}_k = \mathbf{m}_k^{new}$  and repeat. See MacKay (2002) for more discussion.



# K-MEANS ALGORITHM

Start with data  $\mathcal{X} = \{\mathbf{x}^1, \dots, \mathbf{x}^N\}$  and  $K$  starting ‘means’,  $\{\mathbf{m}_1, \dots, \mathbf{m}_K\}$ .

**Assignment:** For each data point,  $\mathbf{x}^n$ , set  $\hat{k}_n = \arg \min_k d(\mathbf{x}^n, \mathbf{m}_k)$ . Set  $\rho_i^n = \delta_i^{\hat{k}_n}$  (Hard responsibility).

**Update:** Let  $N_k = \sum_{n=1}^N \rho_k^n$  and

$$\mathbf{m}_k^{new} = \frac{\sum_{n=1}^N \rho_k^n \mathbf{x}^{(n)}}{N_k}.$$

If  $d(\mathbf{m}_k, \mathbf{m}_k^{new})$  is small for every  $k$ , stop. Otherwise set  $\mathbf{m}_k = \mathbf{m}_k^{new}$  and repeat. See MacKay (2002) for more discussion.

# K-MEANS ALGORITHM

Start with data  $\mathcal{X} = \{\mathbf{x}^1, \dots, \mathbf{x}^N\}$  and  $K$  starting ‘means’,  $\{\mathbf{m}_1, \dots, \mathbf{m}_K\}$ .

**Assignment:** For each data point,  $\mathbf{x}^n$ , set  $\hat{k}_n = \arg \min_k d(\mathbf{x}^n, \mathbf{m}_k)$ . Set  $\rho_i^n = \delta_i^{\hat{k}_n}$  (Hard responsibility).

**Update:** Let  $N_k = \sum_{n=1}^N \rho_k^n$  and

$$\mathbf{m}_k^{new} = \frac{\sum_{n=1}^N \rho_k^n \mathbf{x}^{(n)}}{N_k}.$$

If  $d(\mathbf{m}_k, \mathbf{m}_k^{new})$  is small for every  $k$ , stop. Otherwise set  $\mathbf{m}_k = \mathbf{m}_k^{new}$  and repeat. See MacKay (2002) for more discussion.

# $K$ -MEANS ALGORITHM

Start with data  $\mathcal{X} = \{\mathbf{x}^1, \dots, \mathbf{x}^N\}$  and  $K$  starting ‘means’,  $\{\mathbf{m}_1, \dots, \mathbf{m}_K\}$ .

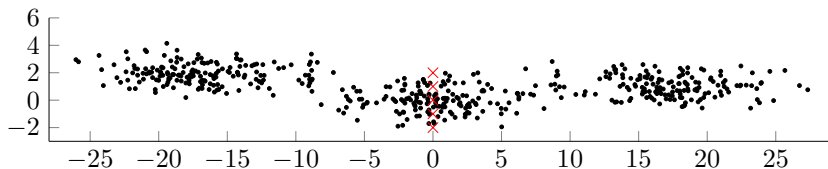
**Assignment:** For each data point,  $\mathbf{x}^n$ , set  $\hat{k}_n = \arg \min_k d(\mathbf{x}^n, \mathbf{m}_k)$ . Set  $\rho_i^n = \delta_i^{\hat{k}_n}$  (Hard responsibility).

**Update:** Let  $N_k = \sum_{n=1}^N \rho_k^n$  and

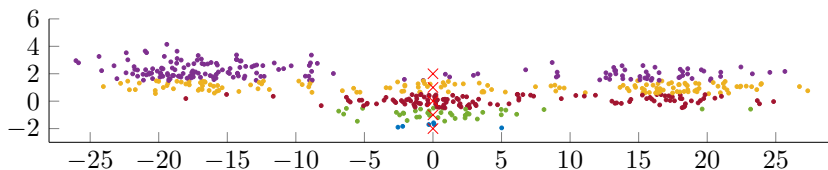
$$\mathbf{m}_k^{new} = \frac{\sum_{n=1}^N \rho_k^n \mathbf{x}^{(n)}}{N_k}.$$

If  $d(\mathbf{m}_k, \mathbf{m}_k^{new})$  is small for every  $k$ , stop. Otherwise set  $\mathbf{m}_k = \mathbf{m}_k^{new}$  and **repeat**. See MacKay (2002) for more discussion.

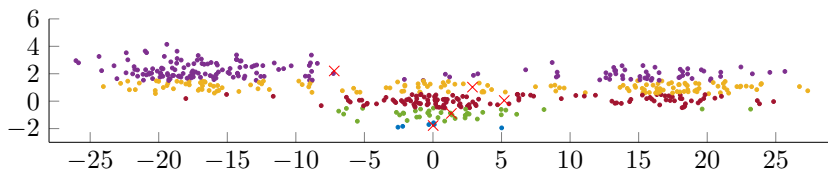
# $K$ -MEANS EXAMPLE



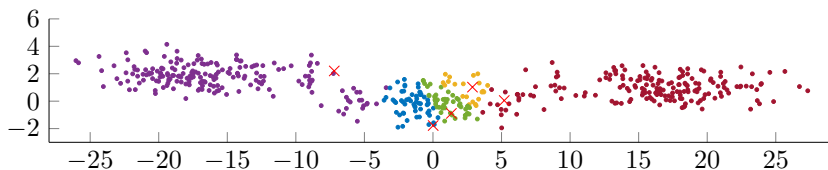
# $K$ -MEANS EXAMPLE



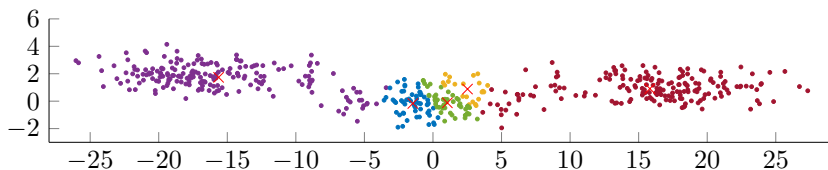
# $K$ -MEANS EXAMPLE



# $K$ -MEANS EXAMPLE

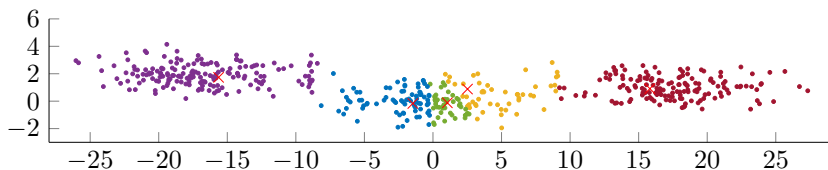


# $K$ -MEANS EXAMPLE

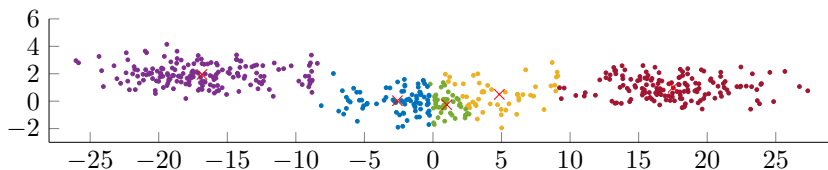




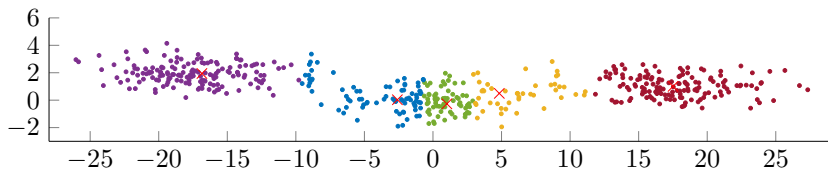
# $K$ -MEANS EXAMPLE



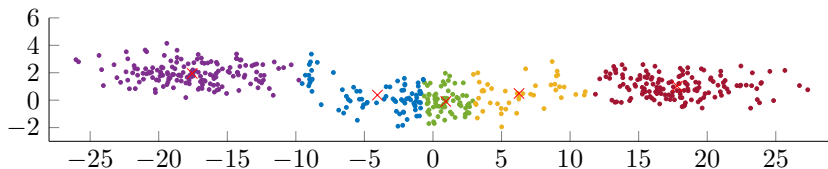
# $K$ -MEANS EXAMPLE



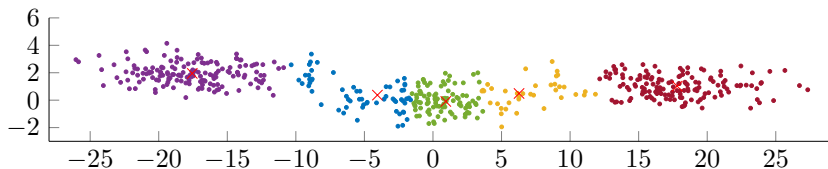
# $K$ -MEANS EXAMPLE



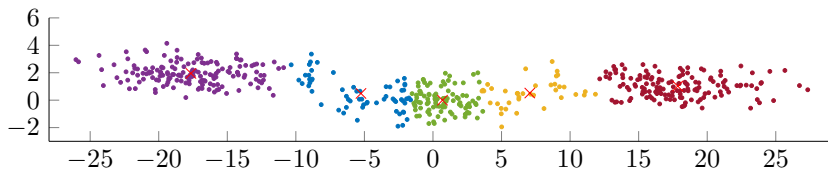
# $K$ -MEANS EXAMPLE



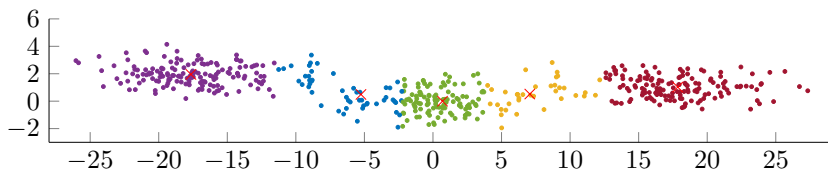
# $K$ -MEANS EXAMPLE



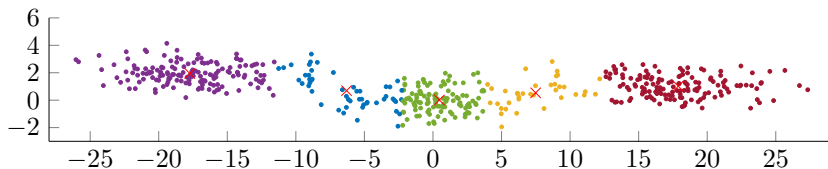
# K-MEANS EXAMPLE



# $K$ -MEANS EXAMPLE

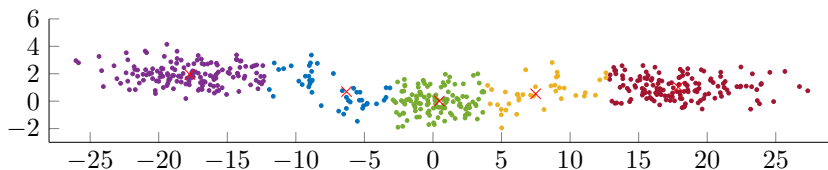


# $K$ -MEANS EXAMPLE

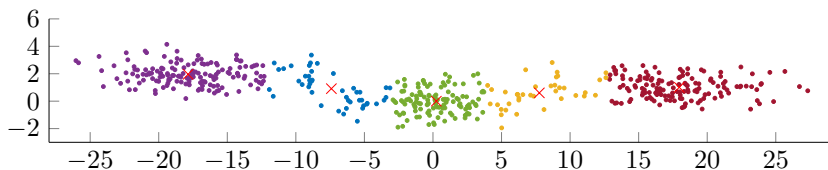




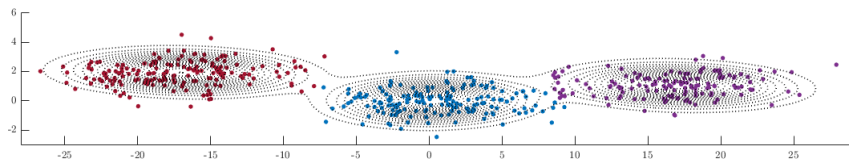
# $K$ -MEANS EXAMPLE



# $K$ -MEANS EXAMPLE

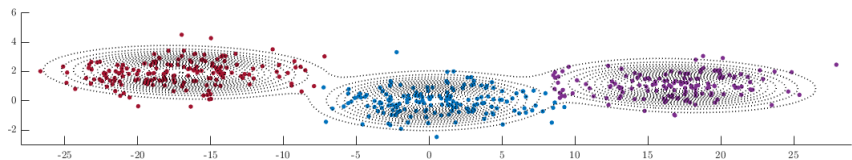


# EM ALGORITHM ON GAUSSIAN MIXTURE

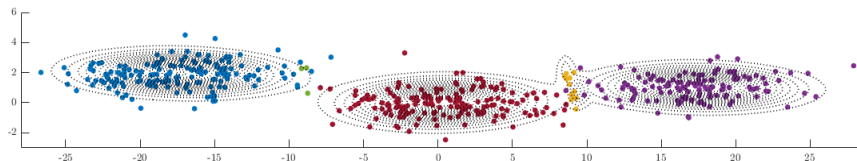


(a)  $AIC \approx 5030$ ,  $BIC \approx 5102$

# EM ALGORITHM ON GAUSSIAN MIXTURE

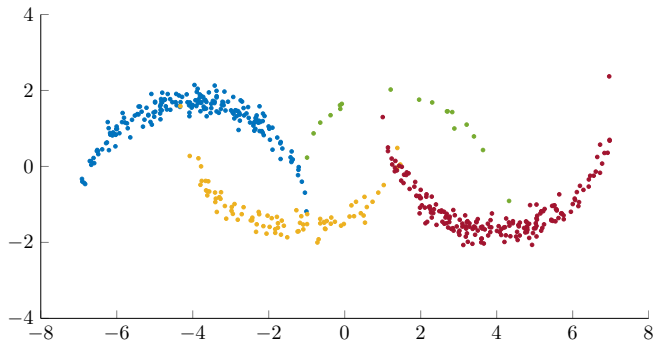


(a)  $AIC \approx 5030$ ,  $BIC \approx 5102$

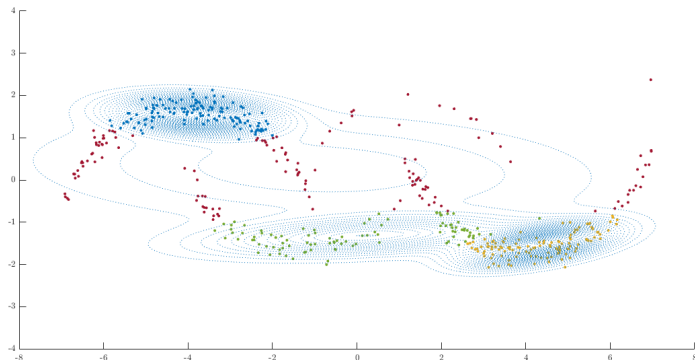


(b)  $AIC \approx 4994$ ,  $BIC \approx 5116$

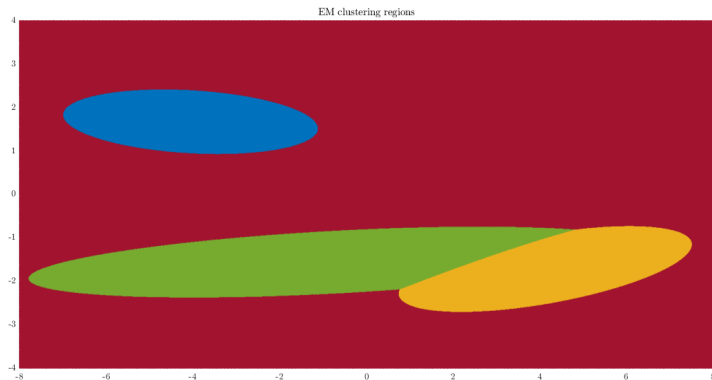
# EM ALGORITHM ON NON-GAUSSIAN DATA



# EM ALGORITHM ON NON-GAUSSIAN DATA



# EM ALGORITHM ON NON-GAUSSIAN DATA



# OVERVIEW

Clustering and Classification

Dynamic Responsibility

Responsible Softmax

Basic Experiments



# RESPONSIBILITY REQUIREMENTS

## Necessary items

- ▶ Data  $\{\mathcal{X}, \mathcal{T}\} = (\mathbf{x}^n, t^n) \ n = 1, \dots, N$
- ▶ Distributions  $f_k(\mathbf{x}, \boldsymbol{\theta}_k) \ k = 1, \dots, K$
- ▶ Parameter matrix  $F = (f_i(\mathbf{x}^j, \boldsymbol{\theta}_i))_i^j = (F_i^j)$
- ▶ Mixture probabilities  $\boldsymbol{\pi}_0 = (\pi_1, \dots, \pi_K) \in S_K$

## Definition (Probability Simplex)

$$S_K := \left\{ \{\pi_k\}_{k=1}^K : 0 \leq \pi_k \leq 1; \sum_{k=1}^K \pi_k = 1 \right\}.$$

# RESPONSIBILITY REQUIREMENTS

## Necessary items

- ▶ Data  $\{\mathcal{X}, \mathcal{T}\} = (\mathbf{x}^n, t^n) \ n = 1, \dots, N$
- ▶ Distributions  $f_k(\mathbf{x}, \boldsymbol{\theta}_k) \ k = 1, \dots, K$
- ▶ Parameter matrix  $F = (f_i(\mathbf{x}^j, \boldsymbol{\theta}_i))_i^j = (F_i^j)$
- ▶ Mixture probabilities  $\boldsymbol{\pi}_0 = (\pi_1, \dots, \pi_K) \in S_K$

## Definition (Probability Simplex)

$$S_K := \left\{ \{\pi_k\}_{k=1}^K : 0 \leq \pi_k \leq 1; \sum_{k=1}^K \pi_k = 1 \right\}.$$

# BAYES' RULE ESTIMATION OF MIXTURE PROBABILITIES

$$\begin{aligned} P(t^n = k | \mathbf{x}^n, \Theta) &= \frac{P(\mathbf{x}^n | t^n = k, \Theta) P(t^n = k | \Theta)}{P(\mathbf{x}^n | \Theta)} \\ &= \frac{f_k(\mathbf{x}^n, \theta_k) \pi_k}{\sum_i \pi_i f_i(\mathbf{x}^n, \theta_i)} \end{aligned}$$

# RESPONSIBILITY

Start with rational maps

$$r_i(\pi) = \frac{1}{N} \sum_n \frac{\pi_i f_i(\mathbf{x}^n, \boldsymbol{\theta}_k)}{\sum_k \pi_k f_k(\mathbf{x}^n, \boldsymbol{\theta}_k)} \quad i = 1, \dots, K$$

Definition (Responsibility Map)

$$R : S_K \rightarrow S_K : R(\pi_1, \pi_2, \dots, \pi_K) = (r_1(\pi), r_2(\pi), \dots, r_K(\pi)). \quad (2.1)$$

When necessary, write  $R_F(\pi)$  to emphasize dependence on  $K \times N$  parameter matrix  $F$ .

# RESPONSIBILITY

Start with rational maps

$$r_i(\boldsymbol{\pi}) = \frac{1}{N} \sum_n \frac{\pi_i f_i(\mathbf{x}^n, \boldsymbol{\theta}_i)}{\sum_k \pi_k f_k(\mathbf{x}^n, \boldsymbol{\theta}_k)} \quad i = 1, \dots, K$$

Definition (Responsibility Map)

$$R : S_K \rightarrow S_K : R(\pi_1, \pi_2, \dots, \pi_K) = (r_1(\boldsymbol{\pi}), r_2(\boldsymbol{\pi}), \dots, r_K(\boldsymbol{\pi})). \quad (2.1)$$

When necessary, write  $R_F(\boldsymbol{\pi})$  to emphasize dependence on  $K \times N$  parameter matrix  $F$ .

# DYNAMIC RESPONSIBILITY

---

## Algorithm 1 Dynamic Responsibility Algorithm

---

**Require:**  $F$  a  $K \times N$  matrix

**Require:**  $\pi_0, \epsilon$

▷  $\epsilon$  creates halt condition

```
1: procedure ITERATION( $F, \pi_0, \epsilon$ )  
2:    $n \leftarrow 1, \pi_n \leftarrow R_F(\pi_0)$   
3:    $orbit \leftarrow \{\pi_0, \pi_1\}$   
4:   while  $|\pi_n - \pi_{n-1}| > \epsilon |\pi_n|$  do  
5:      $\pi_{n+1} \leftarrow R_F(\pi_n)$   
6:      $orbit \leftarrow \{\pi_0, \dots, \pi_{n+1}\}$   
7:      $n \leftarrow n + 1$   
8:   end while  
9:   return  $orbit$   
10: end procedure
```

▷ at this point  $\pi_{n-1} \approx \hat{\pi}$

# LYAPUNOV FUNCTION

For  $\pi \in \mathbb{R}_+^K$  the positive orthant of  $\mathbb{R}^K$ , let

$$\ell_F(\pi) = \frac{1}{N} \sum_{n=1}^N \log \left( \sum_{k=1}^K \pi_k F_k^n \right)$$

## Lemma

*$-\ell_F(\pi)$  is a Lyapunov function for dynamic responsibility. In other words,*

$$\ell_F(R_F(\pi)) \geq \ell_F(\pi)$$

*With equality if and only if  $R_F(\pi) = \pi$ .*

*Note that if  $F$  has full rank,  $-\ell_F$  is strictly convex.*

# LYAPUNOV FUNCTION

For  $\pi \in \mathbb{R}_+^K$  the positive orthant of  $\mathbb{R}^K$ , let

$$\ell_F(\pi) = \frac{1}{N} \sum_{n=1}^N \log \left( \sum_{k=1}^K \pi_k F_k^n \right)$$

## Lemma

$-\ell_F(\pi)$  is a Lyapunov function for dynamic responsibility. In other words,

$$\ell_F(R_F(\pi)) \geq \ell_F(\pi)$$

With equality if and only if  $R_F(\pi) = \pi$ .

Note that if  $F$  has full rank,  $-\ell_F$  is strictly convex.



# LYAPUNOV FUNCTION

For  $\boldsymbol{\pi} \in \mathbb{R}_+^K$  the positive orthant of  $\mathbb{R}^K$ , let

$$\ell_F(\boldsymbol{\pi}) = \frac{1}{N} \sum_{n=1}^N \log \left( \sum_{k=1}^K \pi_k F_k^n \right)$$

## Lemma

$-\ell_F(\boldsymbol{\pi})$  is a Lyapunov function for dynamic responsibility. In other words,

$$\ell_F(R_F(\boldsymbol{\pi})) \geq \ell_F(\boldsymbol{\pi})$$

With equality if and only if  $R_F(\boldsymbol{\pi}) = \boldsymbol{\pi}$ .

Note that if  $F$  has full rank,  $-\ell_F$  is strictly convex.

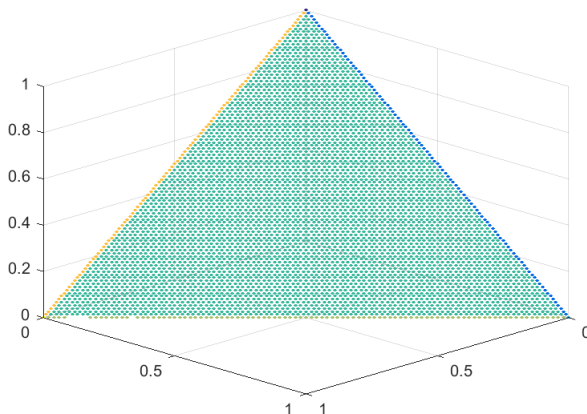
# MAIN THEOREM

Theorem (Convergence of dynamic responsibility )

*If  $F$  has full rank, and  $\pi_0 \in \text{Int } S_K$  then the orbit  $\pi^n = R_F^n(\pi_0)$  converges to  $\hat{\pi}_F$ , the unique maximizing fixed point of  $\ell_F(\pi)$  on  $S_K$ . Moreover,  $\hat{\pi}_F$  depends differentiably on  $F$ .*

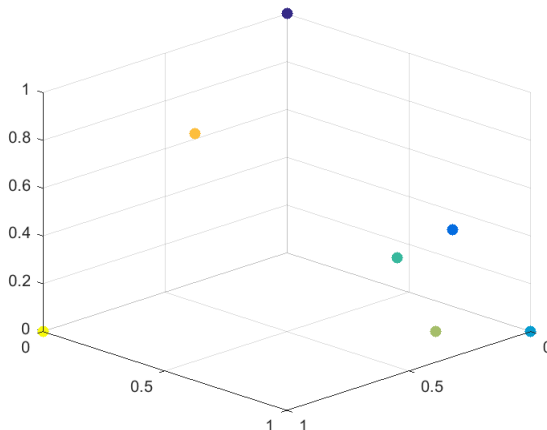
# DYNAMIC RESPONSIBILITY EXAMPLE

If  $F = (F_i^j)$  has linearly independent rows, the interior of  $S_K$  converges to one point.



# DYNAMIC RESPONSIBILITY EXAMPLE

In this case, convergence happens very quickly. (about 5 iterations)



# OVERVIEW

Clustering and Classification

Dynamic Responsibility

Responsible Softmax

Basic Experiments

# CALCULATE $F$

If  $F = e^A$  for some  $A = (A_i^j)$  and  $\mu_i = \ln(\pi_i)$ , then

$$r_i(\boldsymbol{\pi}) = \frac{1}{N} \sum_n \frac{\pi_i F_i^n}{\sum_k \pi_k F_k^n} = \frac{1}{N} \sum_n \frac{\exp(A_i^n + \mu_i)}{\sum_k \exp(A_k^n + \mu_k)}$$

The softmax function is given by the Gibbs Distribution

$$\sigma_i(\boldsymbol{x}) = \frac{\exp(x_i)}{\sum_k \exp(x_k)}.$$

This establishes a connection with modern neural networks.

# CALCULATE $F$

If  $F = e^A$  for some  $A = (A_i^j)$  and  $\mu_i = \ln(\pi_i)$ , then

$$r_i(\boldsymbol{\pi}) = \frac{1}{N} \sum_n \frac{\pi_i F_i^n}{\sum_k \pi_k F_k^n} = \frac{1}{N} \sum_n \frac{\exp(A_i^n + \mu_i)}{\sum_k \exp(A_k^n + \mu_k)}$$

The softmax function is given by the Gibbs Distribution

$$\sigma_i(\boldsymbol{x}) = \frac{\exp(x_i)}{\sum_k \exp(x_k)}.$$

This establishes a connection with modern neural networks.

# CALCULATE $F$

If  $F = e^A$  for some  $A = (A_i^j)$  and  $\mu_i = \ln(\pi_i)$ , then

$$r_i(\boldsymbol{\pi}) = \frac{1}{N} \sum_n \frac{\pi_i F_i^n}{\sum_k \pi_k F_k^n} = \frac{1}{N} \sum_n \frac{\exp(A_i^n + \mu_i)}{\sum_k \exp(A_k^n + \mu_k)}$$

The softmax function is given by the Gibbs Distribution

$$\sigma_i(\boldsymbol{x}) = \frac{\exp(x_i)}{\sum_k \exp(x_k)}.$$

This establishes a connection with modern neural networks.



# CALCULATE $F$

If  $F = e^A$  for some  $A = (A_i^j)$  and  $\mu_i = \ln(\pi_i)$ , then

$$r_i(\boldsymbol{\pi}) = \frac{1}{N} \sum_n \frac{\pi_i F_i^n}{\sum_k \pi_k F_k^n} = \frac{1}{N} \sum_n \frac{\exp(A_i^n + \mu_i)}{\sum_k \exp(A_k^n + \mu_k)}$$

The softmax function is given by the Gibbs Distribution

$$\sigma_i(\boldsymbol{x}) = \frac{\exp(x_i)}{\sum_k \exp(x_k)}.$$

This establishes a connection with modern neural networks.

# NEURAL NETWORK OUTPUT

Neural networks take in data, and output guesses of cluster assignments.

$$F = (F_i^j); \quad \pi^n = R_F^n(\pi_0); \quad \pi^n \rightarrow \hat{\pi} \text{ as } n \rightarrow \infty$$

$$Y(F, \hat{\pi}) = \left( \frac{\hat{\pi}_i F_i^j}{\sum_{k=1}^K \hat{\pi}_k F_k^j} \right)_{i=1, \dots, K}^{j=1, \dots, N}$$

The entry  $Y_i^j$  represents the probability that  $x^j$  comes from cluster  $i$ .

For some  $F$ , it may be that  $\hat{\pi}_F \in \partial S_K$ . To prevent this, stop at some finite  $n = C < \infty$  and use  $Y(F, \pi^C)$  as the output. See Neal and Hinton (1998) for inspiration.

# NEURAL NETWORK OUTPUT

Neural networks take in data, and output guesses of cluster assignments.

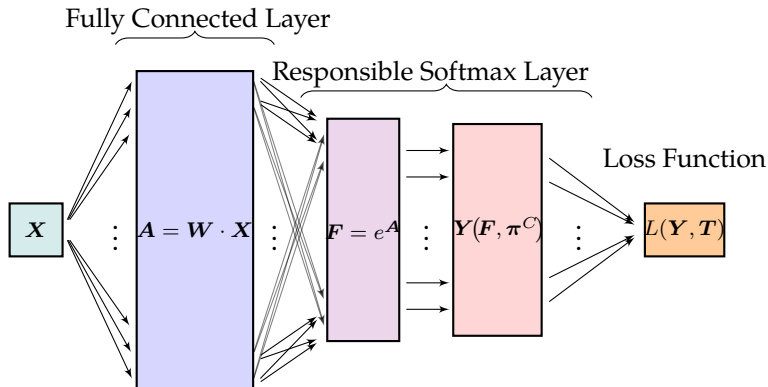
$$F = (F_i^j); \quad \pi^n = R_F^n(\pi_0); \quad \pi^n \rightarrow \hat{\pi} \text{ as } n \rightarrow \infty$$

$$Y(F, \hat{\pi}) = \left( \frac{\hat{\pi}_i F_i^j}{\sum_{k=1}^K \hat{\pi}_k F_k^j} \right)_{i=1, \dots, K}^{j=1, \dots, N}$$

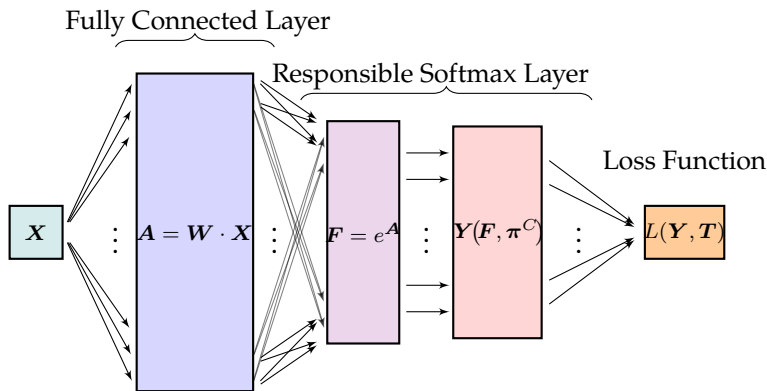
The entry  $Y_i^j$  represents the probability that  $x^j$  comes from cluster  $i$ .

For some  $F$ , it may be that  $\hat{\pi}_F \in \partial S_K$ . To prevent this, stop at some finite  $n = C < \infty$  and use  $Y(F, \pi^C)$  as the output. See Neal and Hinton (1998) for inspiration.

# LAYER DIAGRAM



## LAYER DIAGRAM



$$L(Y, T) = - \sum_n \sum_k T_k^n \log(Y_k^n)$$

# BACKPROPAGATION

The goal is to use gradient descent to learn parameters of the network.

**Option 1:** Automatic differentiation

**Option 2:** Direct calculation

$$\begin{aligned} D\hat{\pi}_F &= D_{\pi}R \cdot D\hat{\pi}_F + D_F R \\ D\hat{\pi}_F &= (I - D_{\pi}R)^{-1} \cdot D_F R \end{aligned} \tag{3.1}$$

In practice, equation (3.1) is too much. An approximation may be used instead.

# BACKPROPAGATION

The goal is to use gradient descent to learn parameters of the network.

**Option 1:** Automatic differentiation

**Option 2:** Direct calculation

$$\begin{aligned} D\hat{\pi}_F &= D_{\pi}R \cdot D\hat{\pi}_F + D_F R \\ D\hat{\pi}_F &= (I - D_{\pi}R)^{-1} \cdot D_F R \end{aligned} \tag{3.1}$$

In practice, equation (3.1) is too much. An approximation may be used instead.  $(I - D_{\pi}R)^{-1} \approx I + DR + DR^2 + \dots + DR^C$

# SETTING THE HYPERPARAMETER $C$

Let  $a_n = d(\pi_{n+1}, \pi_n)$ .

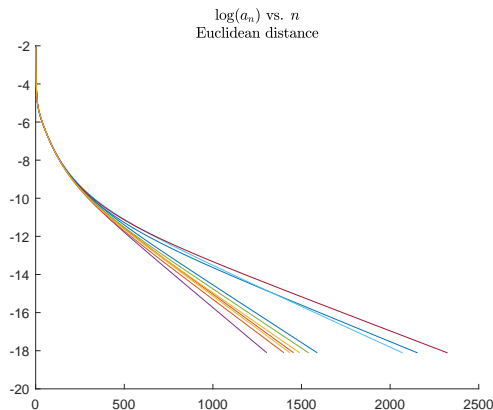


Figure: Plot of  $\log(a_n)$  for several  $F$ . Each curve represents a different parameter matrix  $F$ .



# OVERVIEW

Clustering and Classification

Dynamic Responsibility

Responsible Softmax

Basic Experiments

# EXPERIMENTS WITH GMM

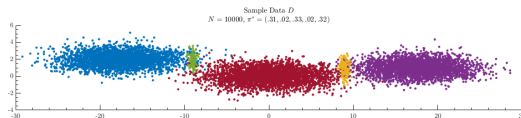
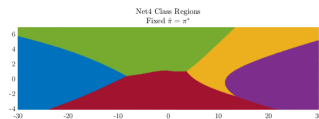
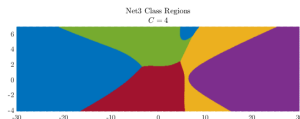
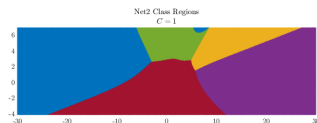
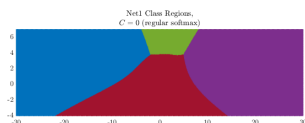


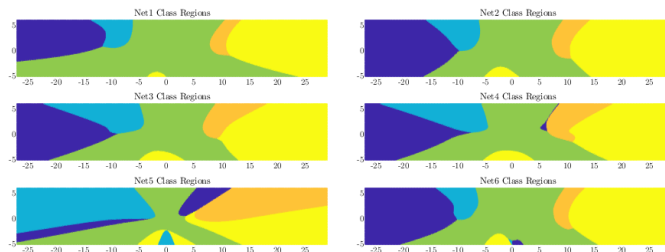
Figure: A sample of data generated from a GMM to test the responsibility softmax layer.

# EXPERIMENTS WITH GMM



Net	Classification layer
Net #1	Softmax
Net #2	Responsibility Softmax; $C = 1$
Net #3	Responsibility Softmax; $C = 4$
Net #4	Fixed Weight Softmax

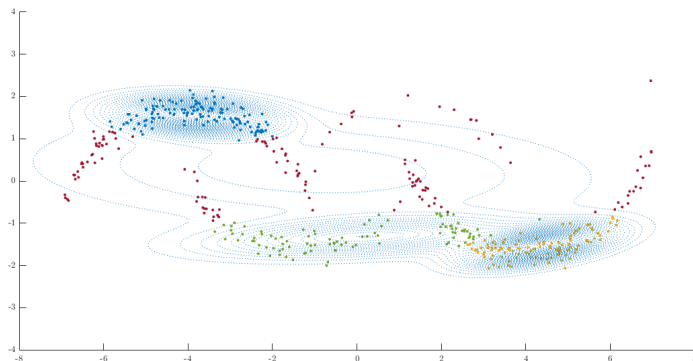
# EXPERIMENTS WITH GMM



Net	Classification layer
Net #1	Softmax
Net #2	Responsibility Softmax; $C = 1$
Net #3	Responsibility Softmax; $C = 4$
Net #4	Responsibility Softmax; $C = 8$
Net #5	Responsibility Softmax; $C = 16$
Net #6	Fixed Weight Softmax

# NON-GAUSSIAN DATA SET

Recall the performance of the EM algorithm on Crescent data



# NON-GAUSSIAN DATA SET

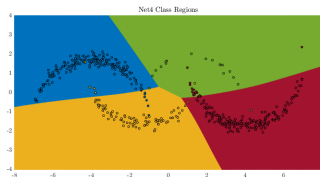
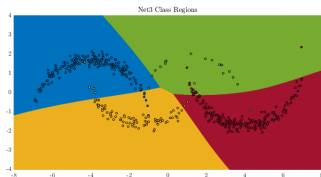
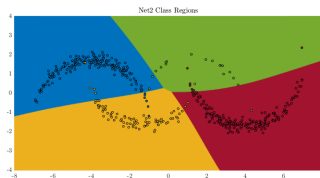
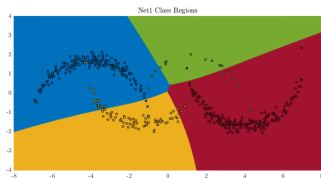


Figure: Classification regions for neural nets trained on crescent data. Hyperparameters are as in GMM example.

# NON-GAUSSIAN DATA SET

**Net 1 Confusion**

1	195	3		
2	15	64		8
3	3		10	5
4			1	196
	1	2	3	4

True Class

Predicted Class

**Net 2 Confusion**

1	192	6		
2	4	78	1	4
3	2		15	1
4			6	191
	1	2	3	4

True Class

Predicted Class

**Net 3 Confusion**

1	192	6		
2	4	78	3	2
3	2		15	1
4			11	186
	1	2	3	4

True Class

Predicted Class

**Net 4 Confusion**

1	192	6		
2	4	79	3	1
3	1		16	1
4			13	184
	1	2	3	4

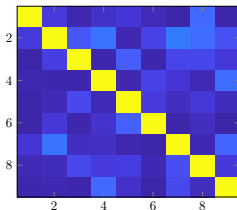
True Class

Predicted Class

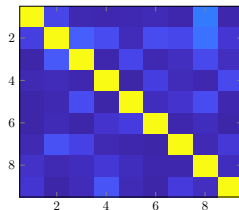
Net	Classification layer
Net #1	Softmax
Net #2	Responsibility Softmax $C = 1$
Net #3	Responsibility Softmax $C = 4$
Net #4	Fixed Weight Softmax

# EXPERIMENTS WITH MNIST

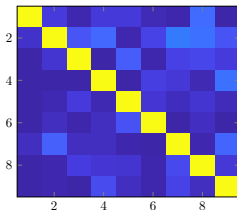
MNIST net 1 Confusion



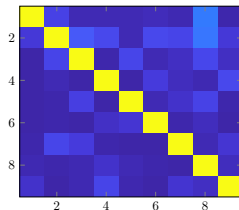
MNIST net 2 Confusion



MNIST net 3 Confusion



MNIST net 4 Confusion





# CONCLUSIONS

We have shown that:

- ▶ **Dynamic responsibility** has nice convergence properties; converges to a MLE.
- ▶ The **responsibility softmax** layer uses dynamic responsibility and gives cluster responsibilities.
- ▶ Using a responsibility softmax layer gives better results when working with imbalanced data. It also works when we do not have distributions for the mixture populations.

# FUTURE WORK

Future work:

- ▶ Use responsibility softmax with other neural nets, LSTM, VAE, Deductron etc.
- ▶ Use responsibility softmax with nonparametric models (e.g. Gaussian processes).
- ▶ Obtain constructive bounds on convergence rates.
- ▶ Explore the relationship between hessian of  $\ell_F$  and Fisher Information matrix.

# REFERENCES

- Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. Information science and statistics. Springer, 1st ed. 2006. corr. 2nd printing edition.
- Deisenroth, M. P., Faisal, A. A., and Ong, C. S. (2020). *Mathematics for Machine Learning*. Cambridge University Press.
- MacKay, D. J. C. (2002). *Information Theory, Inference & Learning Algorithms*. Cambridge University Press, New York, NY, USA.
- Neal, R. M. and Hinton, G. E. (1998). *A View of the EM Algorithm that Justifies Incremental, Sparse, and other Variants*, pages 355–368. Springer Netherlands, Dordrecht.

# EM ALGORITHM FOR GMM

$$\mathcal{X} = \{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^N\}$$

$$f_k(\mathbf{x}) \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k), \quad k = 1, \dots, K$$

$$p(t^n = k) = \pi_k, \quad \sum_k \pi_k = 1$$

# EM ALGORITHM FOR GMM

1. **Expectation** step: Set

$$\rho_k^n = \frac{\pi_k f_k(\mathbf{x}^{(n)})}{\sum_{j=1}^K \pi_j f_j(\mathbf{x}^{(n)})}$$

2. **Maximization** step: Set

$$N_k = \sum_{n=1}^N \rho_k^n, \quad \boldsymbol{\mu}_k^{new} = \frac{1}{N_k} \sum_{n=1}^N \rho_k^n \mathbf{x}^{(n)}$$
$$\pi_k^{new} = \frac{N_k}{N}, \quad \boldsymbol{\Sigma}_k^{new} = \frac{1}{N_k} \sum_{n=1}^N \rho_k^n (\mathbf{x}^{(n)} - \boldsymbol{\mu}_k^{new})(\mathbf{x}^{(n)} - \boldsymbol{\mu}_k^{new})^\top$$

3. Repeat steps 1 and 2 until convergence.

See Bishop (2006) for more details.

# PROOF OF LYAPUNOV LEMMA

## Lemma

*The map  $R_F(\boldsymbol{\pi})$  as defined in equation (2.1) satisfies*

$$R_F(\boldsymbol{\pi}) = \left( \pi_i \cdot \frac{\partial \ell_F}{\partial \pi_i} \bigg|_{\boldsymbol{\pi}} \right)_{1 \leq i \leq K}$$

# PROOF OF LYAPUNOV LEMMA

$$\begin{aligned}
 \ell_F(R_F(\boldsymbol{\pi})) - \ell_F(\boldsymbol{\pi}) &= \frac{1}{N} \sum_{n=1}^N \log \left\{ \frac{\sum_{i=1}^K \pi_i F_i^n \frac{\partial \ell}{\partial \pi_i}}{\sum_{k=1}^K \pi_k F_k^n} \right\} \\
 &\geq \sum_{n=1}^N \sum_{i=1}^K \frac{1}{N} \frac{\pi_i F_i^n}{\sum_{k=1}^K \pi_k f_{kn}} \log \left( \frac{\partial \ell}{\partial \pi_i} \right) \\
 &= \sum_{i=1}^K \sum_{n=1}^N \frac{1}{N} \frac{\pi_i F_i^n}{\sum_{k=1}^K \pi_k f_{kn}} \log \left( \frac{\partial \ell}{\partial \pi_i} \right) \\
 &= \sum_{i=1}^K r_i(\boldsymbol{\pi}) \log \left( \frac{r_i(\boldsymbol{\pi})}{\pi_i} \right) \geq 0
 \end{aligned}$$

# PROOF OF LYAPUNOV LEMMA

$$\begin{aligned}
 \ell_F(R_F(\boldsymbol{\pi})) - \ell_F(\boldsymbol{\pi}) &= \frac{1}{N} \sum_{n=1}^N \log \left\{ \frac{\sum_{i=1}^K \pi_i F_i^n \frac{\partial \ell}{\partial \pi_i}}{\sum_{k=1}^K \pi_k F_k^n} \right\} \\
 &\geq \sum_{n=1}^N \sum_{i=1}^K \frac{1}{N} \frac{\pi_i F_i^n}{\sum_{k=1}^K \pi_k f_{kn}} \log \left( \frac{\partial \ell}{\partial \pi_i} \right) \\
 &= \sum_{i=1}^K \sum_{n=1}^N \frac{1}{N} \frac{\pi_i F_i^n}{\sum_{k=1}^K \pi_k f_{kn}} \log \left( \frac{\partial \ell}{\partial \pi_i} \right) \\
 &= \sum_{i=1}^K r_i(\boldsymbol{\pi}) \log \left( \frac{r_i(\boldsymbol{\pi})}{\pi_i} \right) \geq 0
 \end{aligned}$$



# PROOF OF LYAPUNOV LEMMA

$$\begin{aligned}
 \ell_F(R_F(\boldsymbol{\pi})) - \ell_F(\boldsymbol{\pi}) &= \frac{1}{N} \sum_{n=1}^N \log \left\{ \frac{\sum_{i=1}^K \pi_i F_i^n \frac{\partial \ell}{\partial \pi_i}}{\sum_{k=1}^K \pi_k F_k^n} \right\} \\
 &\geq \sum_{n=1}^N \sum_{i=1}^K \frac{1}{N} \frac{\pi_i F_i^n}{\sum_{k=1}^K \pi_k f_{kn}} \log \left( \frac{\partial \ell}{\partial \pi_i} \right) \\
 &= \sum_{i=1}^K \sum_{n=1}^N \frac{1}{N} \frac{\pi_i F_i^n}{\sum_{k=1}^K \pi_k f_{kn}} \log \left( \frac{\partial \ell}{\partial \pi_i} \right) \\
 &= \sum_{i=1}^K r_i(\boldsymbol{\pi}) \log \left( \frac{r_i(\boldsymbol{\pi})}{\pi_i} \right) \geq 0
 \end{aligned}$$

# PROOF OF LYAPUNOV LEMMA

$$\begin{aligned}
 \ell_F(R_F(\boldsymbol{\pi})) - \ell_F(\boldsymbol{\pi}) &= \frac{1}{N} \sum_{n=1}^N \log \left\{ \frac{\sum_{i=1}^K \pi_i F_i^n \frac{\partial \ell}{\partial \pi_i}}{\sum_{k=1}^K \pi_k F_k^n} \right\} \\
 &\geq \sum_{n=1}^N \sum_{i=1}^K \frac{1}{N} \frac{\pi_i F_i^n}{\sum_{k=1}^K \pi_k f_{kn}} \log \left( \frac{\partial \ell}{\partial \pi_i} \right) \\
 &= \sum_{i=1}^K \sum_{n=1}^N \frac{1}{N} \frac{\pi_i F_i^n}{\sum_{k=1}^K \pi_k f_{kn}} \log \left( \frac{\partial \ell}{\partial \pi_i} \right) \\
 &= \sum_{i=1}^K r_i(\boldsymbol{\pi}) \log \left( \frac{r_i(\boldsymbol{\pi})}{\pi_i} \right) \geq 0
 \end{aligned}$$

## CONFUSION FOR GMM DATA

$30.253 \pm .001$	0.0	$.027 \pm .001$	0.0	0.0
$1.680 \pm .000$	0.0	0.0	0.0	0.0
$.328 \pm .004$	0.0	$32.206 \pm .008$	0.0	$.706 \pm .006$
0.0	0.0	$.033 \pm .001$	0.0	$3.207 \pm .001$
0.0	0.0	$.021 \pm .001$	0.0	$31.539 \pm .001$

(a) Confusion table for GMM Net #1.

Figure: The nets were tested on a set of samples drawn independently from the training set. Values are reported as percentages for clarity. Test data sample size  $N = 2500$  for all runs. Error intervals are 95% confidence standard error. An entry of 0.0 indicates that all values were zero to 3 decimal places.

## CONFUSION FOR GMM DATA

30.165 $\pm$ .004	.101 $\pm$ .004	.014 $\pm$ .001	0.0	0.0
1.616 $\pm$ .003	.063 $\pm$ .003	0.0	0.0	0.0
.398 $\pm$ .006	.114 $\pm$ .003	31.739 $\pm$ .010	.330 $\pm$ .008	.659 $\pm$ .009
0.0	0.0	.031 $\pm$ .001	.333 $\pm$ .012	2.875 $\pm$ .012
0.0	0.0	.012 $\pm$ .000	.082 $\pm$ .004	31.466 $\pm$ .004

(a) Confusion table for GMM Net #2.

Figure: The nets were tested on a set of samples drawn independently from the training set. Values are reported as percentages for clarity. Test data sample size  $N = 2500$  for all runs. Error intervals are 95% confidence standard error. An entry of 0.0 indicates that all values were zero to 3 decimal places.

# CONFUSION FOR GMM DATA

$29.897 \pm .010$	$.374 \pm .010$	$.009 \pm .001$	$.000 \pm .001$	0.0
$1.273 \pm .011$	$.406 \pm .011$	$.001 \pm .001$	0.0	0.0
$.658 \pm .016$	$.595 \pm .017$	$29.916 \pm .036$	$1.329 \pm .031$	$.743 \pm .018$
0.0	$.000 \pm .001$	$.013 \pm .001$	$1.221 \pm .027$	$2.006 \pm .027$
0.0	0.0	0.0	$.340 \pm .009$	$31.220 \pm .009$

(a) Confusion table for GMM Net #3.

Figure: The nets were tested on a set of samples drawn independently from the training set. Values are reported as percentages for clarity. Test data sample size  $N = 2500$  for all runs. Error intervals are 95% confidence standard error. An entry of 0.0 indicates that all values were zero to 3 decimal places.

## CONFUSION FOR GMM DATA

26.842 $\pm$ .035	3.438 $\pm$ .035	0.0	0.0	0.0
.044 $\pm$ .006	1.636 $\pm$ .006	0.0	0.0	0.0
.075 $\pm$ .003	1.841 $\pm$ .024	28.737 $\pm$ .037	2.540 $\pm$ .025	.047 $\pm$ .002
0.0	0.0	.027 $\pm$ .001	3.122 $\pm$ .004	.092 $\pm$ .004
0.0	0.0	0.0	2.463 $\pm$ .023	29.097 $\pm$ .023

(a) Confusion table for GMM Net #4.

Figure: The nets were tested on a set of samples drawn independently from the training set. Values are reported as percentages for clarity. Test data sample size  $N = 2500$  for all runs. Error intervals are 95% confidence standard error. An entry of 0.0 indicates that all values were zero to 3 decimal places.

## PER CLASS PRECISION AND RECALL FOR GMM DATA

GMM Net 1		
Class	Precision	Recall
1	0.936	0.999
2	0.000	0.000
3	0.998	0.966
4	0.000	0.000
5	0.979	0.999

(a) Precision and Recall table for GMM Net #1.

Figure: This table shows per class precision and recall for GMM nets trained and tested on the same data as in table 5

## PER CLASS PRECISION AND RECALL FOR GMM DATA

GMM Net 2		
Class	Precision	Recall
1	0.935	0.997
2	0.209	0.027
3	0.998	0.953
4	0.401	0.090
5	0.898	0.997

(a) Precision and Recall table for GMM Net #2.

Figure: This table shows per class precision and recall for GMM nets trained and tested on the same data as in table 5



## PER CLASS PRECISION AND RECALL FOR GMM DATA

GMM Net 3		
Class	Precision	Recall
1	0.934	0.988
2	0.279	0.194
3	0.999	0.894
4	0.385	0.364
5	0.918	0.989

(a) Precision and Recall table for GMM Net #3.

Figure: This table shows per class precision and recall for GMM nets trained and tested on the same data as in table 5

## PER CLASS PRECISION AND RECALL FOR GMM DATA

GMM Net 4		
Class	Precision	Recall
1	0.988	0.930
2	0.289	0.882
3	0.999	0.868
4	0.380	0.969
5	0.996	0.922

(a) Precision and Recall table for GMM Net #4.

Figure: This table shows per class precision and recall for GMM nets trained and tested on the same data as in table 5