

Efficient Maximal Coding Rate Reduction by Variational Forms

-Summary-

Introduction

- Arising new objective for classification : MCR^2 (Maximizing the coding rate reduction)
- Problem of CE :
 - Learning based on CE will leads to **neural collapse** : when $CE \rightarrow 0$, the representations of each class at penultimate layer collapse to a single point, suppressing within-class variability. (Common symptom on penultimate layer)
- To alleviate this phenomenon, [Yu et al., 2020] suggested **MCR^2 objective** that encourages the latent representation of the entire training set to occupy as much volume as possible, while enforcing each class to occupy as little space as possible.
- Empirically, and theoretically, it was shown that the latent representations will be a low dimensional linear subspace with the subspace orthogonal to each other.

Introduction

- But, calculating MCR^2 is intractable as the # of class k increases (due to $\log \det$). This paper tries to resolve this problem by suggesting variational form of MCR^2 .

- Original MCR^2 objective :

$$\max_{\theta} \Delta R(\mathbf{Z}_{\theta}) \equiv R(\mathbf{Z}_{\theta}) - R_c(\mathbf{Z}_{\theta}, \mathbf{\Pi}) \quad \text{s.t. } \mathbf{Z}_{\theta} \in \mathcal{S}, \quad \text{where}$$

$$R(\mathbf{Z}_{\theta}) = \frac{1}{2} \log \det \left(\mathbf{I} + \alpha \mathbf{Z}_{\theta} \mathbf{Z}_{\theta}^{\top} \right), \quad \text{and}$$

$$R_c(\mathbf{Z}_{\theta}, \mathbf{\Pi}) = \sum_{j=1}^k \frac{\gamma_j}{2} \log \det \left(\mathbf{I} + \alpha_j \mathbf{Z}_{\theta} \text{Diag}(\mathbf{\Pi}_j) \mathbf{Z}_{\theta}^{\top} \right)$$

where $\mathbf{Z} = [f_{\theta}(X_1), \dots, f_{\theta}(X_m)] \in \mathbb{R}^{d \times m}$, $\mathbf{X} = [X_1, \dots, X_m] \in \mathbb{R}^{D \times m}$

and $\mathbf{\Pi} \in \mathbb{R}^{m \times k}$ = class membership matrix

[$\Pi_{i,j} = p(X_i \text{ is in class } j)$]

k = # of classes

m = # of samples

X_i = i th sample

d = feature dimension

D = input dimension

$\mathbf{\Pi}$ = membership matrix

$\mathbf{\Pi}_j$ = j th column of $\mathbf{\Pi}$

Description of some terms :

1. $\alpha = \frac{d}{m\epsilon^2}$

2. $\alpha_j = \frac{d}{\langle \mathbf{1}, \mathbf{\Pi}_j \rangle \epsilon^2}$

3. $\gamma_j = \frac{\langle \mathbf{1}, \mathbf{\Pi}_j \rangle}{m}$

Introduction

- Original MCR^2 objective :

$$\max_{\theta} \Delta R(\mathbf{Z}_{\theta}) \equiv R(\mathbf{Z}_{\theta}) - R_c(\mathbf{Z}_{\theta}, \mathbf{\Pi}) \quad \text{s.t. } \mathbf{Z}_{\theta} \in \mathcal{S}, \quad \text{where}$$

$$R(\mathbf{Z}_{\theta}) = \frac{1}{2} \log \det \left(\mathbf{I} + \alpha \mathbf{Z}_{\theta} \mathbf{Z}_{\theta}^{\top} \right), \quad \text{and}$$

$$R_c(\mathbf{Z}_{\theta}, \mathbf{\Pi}) = \sum_{j=1}^k \frac{\gamma_j}{2} \log \det \left(\mathbf{I} + \alpha_j \mathbf{Z}_{\theta} \text{Diag}(\mathbf{\Pi}_j) \mathbf{Z}_{\theta}^{\top} \right)$$

- $R(\mathbf{Z}_{\theta})$ (expansion term) : captures the dimension (or volume) of the space spanned by \mathbf{Z}_{θ}
(\cong required bits to encode \mathbf{Z}_{θ} with assumption of MN)
- $R_c(\mathbf{Z}_{\theta}, \mathbf{\Pi})$ (compression term) : measures the sum of the dimensions (or volumes) of the data from each class
(\cong required bits to encode \mathbf{Z}_{θ} with assumption of mixture ($\sim \mathbf{\Pi}$) of MN)

k = # of classes

m = # of samples

X_i = i th sample

d = feature dimension

D = input dimension

$\mathbf{\Pi}$ = membership matrix

$\mathbf{\Pi}_j$ = j th column of $\mathbf{\Pi}$

Description of some terms :

1. $\alpha = \frac{d}{m\epsilon^2}$

2. $\alpha_j = \frac{d}{\langle \mathbf{1}, \mathbf{\Pi}_j \rangle \epsilon^2}$

3. $\gamma_j = \frac{\langle \mathbf{1}, \mathbf{\Pi}_j \rangle}{m}$

Introduction

- What does Original MCR^2 objective do?
 - While maximizing the overall volume of the embedded features (1st term), we want to compress the volume of embedded features from each class (2nd term)
 - (In other perspective, it can be seen as minimizing the # of bits to encode Z when Π is known while maximizing the # of bits to encode Z when we don't have mixture assumption)
- But, calculating $R_c(Z_\theta, \Pi)$ involves k computations of $\log \det$, which becomes intractable as # of classes \uparrow .
- They cleverly modify their terms into variational form to avoid this problem.

Variational formulation of MCR^2

- Observe following fact and theorems :

1. Identity equation for PSD matrix M and $c \geq 0$:

$$\log \det(\mathbf{I} + c\mathbf{M}) = \sum_{i=1}^r \log(1 + c\sigma_i(\mathbf{M}))$$

2. Variational form of spectral functions :

Theorem 2.1 (Adapted from [16]) For any matrix \mathbf{X} , let r denote the rank of \mathbf{X} , let $\sigma_i(\mathbf{X})$ denote the i^{th} singular value of \mathbf{X} , and define

$$H(\mathbf{X}) = \sum_{i=1}^r h(\sigma_i(\mathbf{X})).$$

for some function h . If h is a concave, non-decreasing function on $[0, \infty)$ with $h(0) = 0$, then the following holds

$$H(\mathbf{X}) = \min_{\mathbf{U}, \mathbf{V}: \mathbf{U}\mathbf{V}^\top = \mathbf{X}} \sum_i h(\|\mathbf{U}_i\|_2 \|\mathbf{V}_i\|_2),$$

where $(\mathbf{U}_i, \mathbf{V}_i)$ denotes the i^{th} columns of (\mathbf{U}, \mathbf{V}) . Note also that (\mathbf{U}, \mathbf{V}) can have an arbitrary number of columns ($\geq r$) provided $\mathbf{U}\mathbf{V}^\top = \mathbf{X}$.



Proposition 3.1 Let \mathbf{M} be any real positive semi-definite matrix and let $c \geq 0$ be any non-negative scalar. Then the following holds:

$$-\log \det(\mathbf{I} + c\mathbf{M}) = \max_{\mathbf{U}: \mathbf{U}\mathbf{U}^\top = \mathbf{M}} -\sum_i \log(1 + c\|\mathbf{U}_i\|_2^2). \quad (3)$$

Further, if $\bar{\mathbf{U}}\mathbf{S}\bar{\mathbf{U}}^\top = \mathbf{M}$ is a SVD of \mathbf{M} then $\mathbf{U}^* = \bar{\mathbf{U}}\mathbf{S}^{1/2}$ is a solution to the above problem.

Variational formulation of MCR^2

- Using proposition 3.1, we can optimization variable $\{U^{(j)}\}_{j=1}^k$ to replace $Z_\theta \text{Diag}(\Pi_j) Z_\theta^T$:

$$\max_{\theta} \Delta R(Z_\theta) =$$

$$\max_{\theta, \{U^{(j)}\}_{j=1}^k} \frac{1}{2} \log \det \left(I + \alpha \sum_{j=1}^k U^{(j)} (U^{(j)})^\top \right)$$

$$- \sum_{j=1}^k \frac{\gamma_j}{2} \sum_i \log \left(1 + \alpha_j \|U_i^{(j)}\|_2^2 \right)$$

$$\text{s.t. } \forall j, U^{(j)} (U^{(j)})^\top = Z_\theta \text{Diag}(\Pi_j) Z_\theta^\top \text{ and } Z_\theta \in \mathcal{S}.$$

where $U^{(j)} = \Gamma \text{Diag}(A_j)^{\frac{1}{2}}$, $\Gamma \in \mathbb{R}^{d \times q} \cap \mathcal{S}$ (dictionary with unit l_2 normalized columns),
 $A_j \in \mathbb{R}_+^q$ (non-negative encoding vector)

- Due to this reparameterization, $\Gamma \text{Diag}(A_j) \Gamma^T = U^{(j)} U^{(j)T}$ and $\|U_i^{(j)}\|_2^2 = A_{i,j}$ holds

Note :

q is larger than k (usually $q \cong 2k$)

Notation :

\mathcal{S} = matrix sets where their columns are l_2 normalized.

Variational formulation of MCR^2

- In addition to this reparameterization, they add regularization term for this reparameterization :

$$M(\mathbf{Z}_\theta, \Gamma, \mathbf{A}) = \sum_{j=1}^k \frac{1}{\gamma_j} \left\| \mathbf{Z}_\theta \text{Diag}(\Pi_j) \mathbf{Z}_\theta^T - \Gamma \text{Diag}(\mathbf{A}_j) \Gamma^T \right\|_F^2$$

- Their final proposed formulation ($V - MCR^2$) is following :

$$\max_{\theta, \Gamma \in \mathbb{R}^{d \times q} \cap \mathcal{S}, \mathbf{A} \in \mathbb{R}_+^{q \times k}} R^v(\Gamma, \mathbf{A}) - R_c^v(\mathbf{A}) - \frac{\mu}{2m} M(\mathbf{Z}_\theta, \Gamma, \mathbf{A})$$

$$\text{where } R^v(\Gamma, \mathbf{A}) = \frac{1}{2} \log \det \left(\mathbf{I} + \alpha \sum_{j=1}^k \Gamma \text{Diag}(\mathbf{A}_j) \Gamma^T \right),$$

$$R_c^v(\mathbf{A}) = \sum_{j=1}^k \frac{\gamma_j}{2} \sum_{l=1}^q \log(1 + \alpha_j \mathbf{A}_{l,j}),$$

$$M(\mathbf{Z}_\theta, \Gamma, \mathbf{A}) = \sum_{i=1}^k \frac{1}{\gamma_j} \left\| \mathbf{Z}_\theta \text{Diag}(\Pi_j) \mathbf{Z}_\theta^T - \Gamma \text{Diag}(\mathbf{A}_j) \Gamma^T \right\|_F^2$$

Note :

The penalization term M can be interpreted a class-balanced low-rank LASSO

∴ it imposes $\Gamma \text{Diag}(\mathbf{A}_j) \Gamma^T \rightarrow \mathbf{Z}_\theta \text{Diag}(\Pi_j) \mathbf{Z}_\theta^T$
 where $\Gamma \in \mathbb{R}^{d \times q}, \mathbf{A}_j \in \mathbb{R}_+^q, \mathbf{Z}_\theta \in \mathbb{R}^{d \times k}, \Pi_j \in \mathbb{R}_+^k$
 and $q > k$.

Variational formulation of MCR^2

- Now, the complexity of calculating R_C term is changed from $O(k \min\{d^3, m^3\})$ to $O(qk)$.
- *Note : complexity for calculating M is $O(kq^2)$, which is also tractable compared to original R_C terms.*
- Then, how to optimize this terms?
 - 1st step : optimize Γ, A by stable GA (using Lipschitz upper-bounded l_r) using whole objective.
 - 2nd step : optimize θ by naïve GD using only $M(Z_\theta, \Gamma, A)$
 - 3rd step : Guide Γ, A by explicit solution : SVD of $Z_\theta \text{Diag}(\Pi_j) Z_\theta^T$ (called ‘latching’)

Variational formulation of MCR^2

- Overall algorithm :

Algorithm 1 Variational MCR^2 Training

```

1: Input: data  $\mathbf{X}$ , labels  $\mathbf{Y}$ , featurizer  $f_\theta(\cdot)$ ,
   latch-freq, step sizes  $(\nu_\theta, \nu_\Gamma, \nu_{\mathbf{A}})$ 
2: Initialize  $\mathbf{A}, \Gamma \leftarrow \text{latching}(\mathbf{X}, \mathbf{Y}, f_\theta)$ 
3: for iter = 0, 1, ...,  $n - 1$  do
4:   Get  $\mathbf{Z}_\theta = f_\theta(\mathbf{X})$  and membership matrices  $\mathbf{\Pi}$ 
5:   Get  $\ell_{\text{V-MCR}^2}(\mathbf{Z}_\theta, \Gamma, \mathbf{A})$ 
6:   Compute  $L_{\mathbf{A}}, L_\Gamma$ 
7:    $\Gamma \leftarrow \Gamma + \frac{\nu_\Gamma}{L_\Gamma} \nabla_\Gamma \ell_{\text{V-MCR}^2}(\mathbf{Z}_\theta, \Gamma, \mathbf{A})$ 
8:    $\mathbf{A} \leftarrow \mathbf{A} + \frac{\nu_{\mathbf{A}}}{L_{\mathbf{A}}} \nabla_{\mathbf{A}} \ell_{\text{V-MCR}^2}(\mathbf{Z}_\theta, \Gamma, \mathbf{A})$ 
9:   Project  $\mathbf{A} \leftarrow \text{ReLU}(\mathbf{A})$ 
10:  Project  $\Gamma_l \leftarrow \frac{1}{\|\Gamma_l\|_2} \Gamma_l \quad \forall l \in [q]$ 
11:  Recompute  $M(\mathbf{Z}_\theta, \Gamma, \mathbf{A})$ 
12:   $\theta \leftarrow \theta - \nu_\theta \nabla_\theta (M(\mathbf{Z}_\theta, \Gamma, \mathbf{A}))$ 
13:  if iter mod latch-freq = 0 then
14:     $\mathbf{A}, \Gamma \leftarrow \text{latching}(\mathbf{X}, \mathbf{Y}, f_\theta)$ 
15:  end if
16: end for
17: return  $f_\theta$ 

```

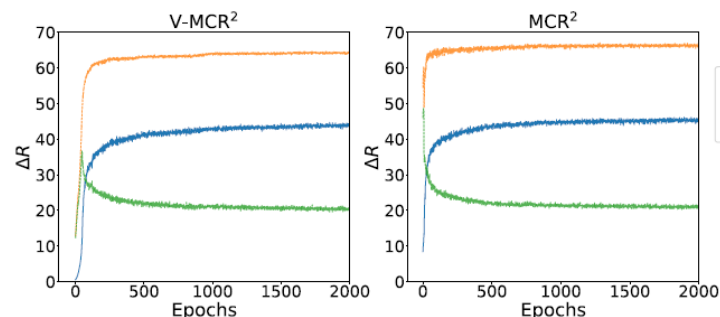
Algorithm 2 Latching

Input: data \mathbf{X} , labels \mathbf{Y} , featurizer $f_\theta(\cdot)$
 Get $\mathbf{Z}_\theta = f_\theta(\mathbf{X}) \in \mathbb{R}^{d \times m}$ and membership $\mathbf{\Pi} \in \mathbb{R}^{m \times k}$
 $\mathbf{A} \leftarrow \mathbf{0} \in \mathbb{R}^{q \times k}$ (assume q is divisible by k)
 $\Gamma \leftarrow \mathbf{0} \in \mathbb{R}^{d \times q}$
for $j = 1, \dots, k$ **do**
 Get $\mathbf{U} \text{Diag}(\boldsymbol{\sigma}) \mathbf{V}^\top = \text{SVD}(\mathbf{Z}_\theta \text{Diag}(\mathbf{\Pi}_j) \mathbf{Z}_\theta^\top)$
 $s \leftarrow q/k$
 $\Gamma[:, (j-1)*s : j*s] = \mathbf{U}[:, 0 : s]$ % python indexing
 $\mathbf{A}[(j-1)*s : j*s, j] = \boldsymbol{\sigma}[0 : s]$ % python indexing
end for
return \mathbf{A}, Γ

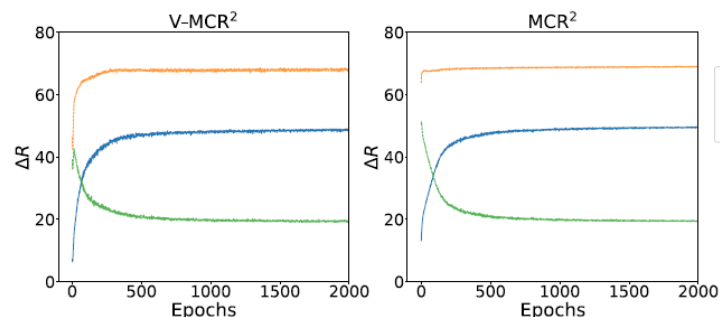
1. Recall that optimal $\mathbf{U}^{(j)} = \text{SVD}(\mathbf{Z}_\theta \text{Diag}(\mathbf{\Pi}_j) \mathbf{Z}_\theta)$ when $q = k$
2. When the latching is applied, we get s^{th} order PC for each $\mathbf{U}^{(j)}$:
 Especially, $\mathbf{U}_{l+j*s}^{(j)} = \sigma_l \mathbf{U}_l$ for each $j \in [k]$ and $l \in [s]$ (o.w 0 vector)
 where $\mathbf{U}_l = l^{\text{th}}$ column of \mathbf{U} (from SVD of $\mathbf{Z}_\theta \text{Diag}(\mathbf{\Pi}_j) \mathbf{Z}_\theta^\top$)
3. According to paper, it helps to improve convergence (although it is quite computationally expensive $\sim O(kd^3)$)

Experiments of $V - MCR^2$

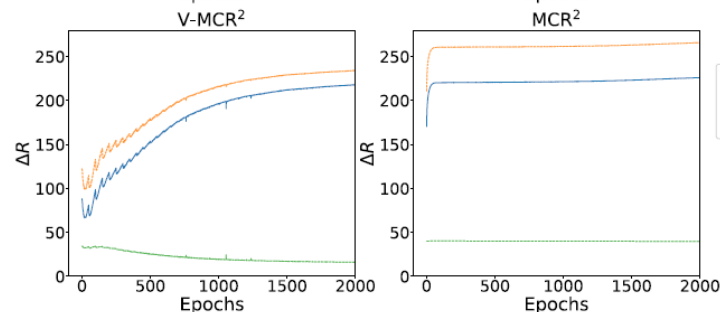
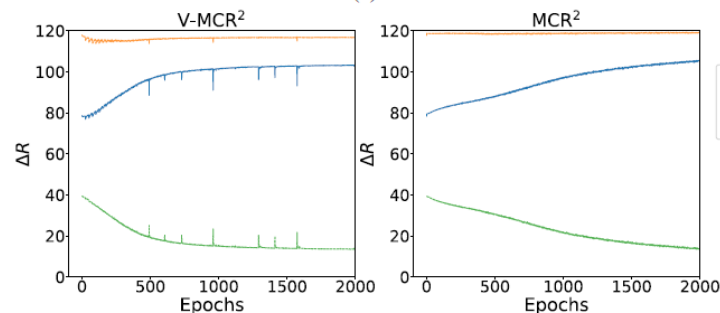
- Then, ΔR of $V - MCR^2$ is similar for original MCR^2 when k and d is small.
- However, it requires some time to have similar ΔR when k and d is large.



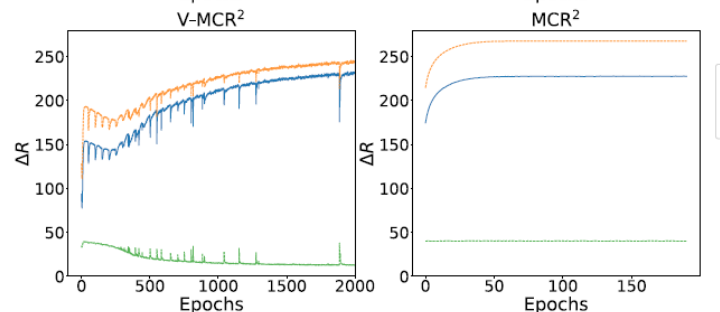
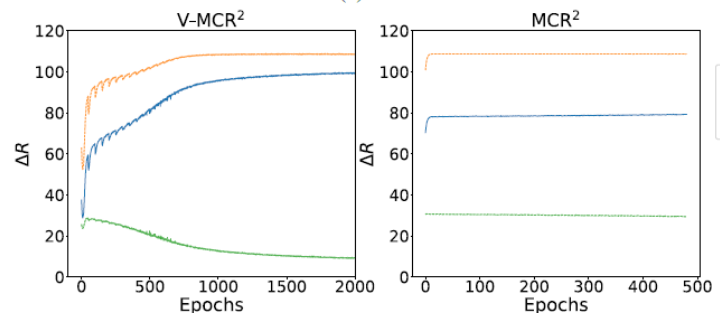
(a) MNIST



(c) CIFAR-10



(b) CIFAR-100 $d = 100$ (top), $d = 500$ (bottom)

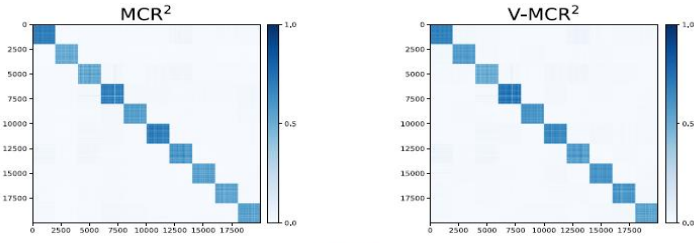


(d) Tiny ImageNet $d = 200$ (top), $d = 500$ (bottom)

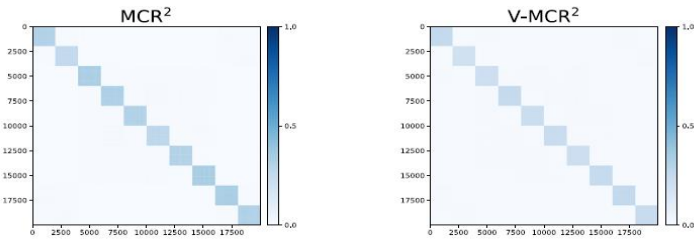
Convergence of training ΔR

Experiments of $V - MCR^2$

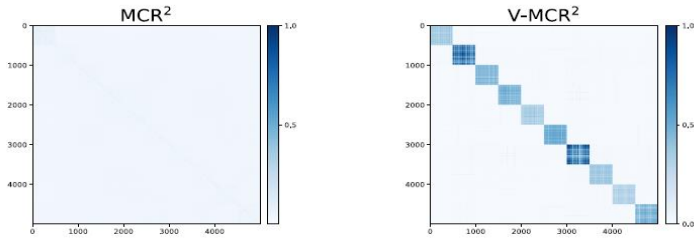
- Surprisingly $V - MCR^2$ shows better representation learning compared to original MCR^2 :



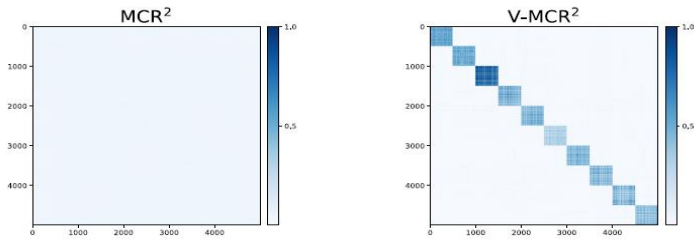
(a) MNIST



(b) CIFAR-10



(c) CIFAR-100



(d) Tiny ImageNet

Dataset	Objective	Training ΔR	Test Accuracy
MNIST	MCR^2	44.6429	0.9785
	V- MCR^2	44.2117	0.9788
	CE	-	0.9738
CIFAR-10	MCR^2	49.40	0.8956
	V- MCR^2	48.43	0.8997
	CE	-	0.8665
CIFAR-100	MCR^2	226.0519	0.2421
	V- MCR^2	218.0185	<u>0.5872</u>
	CE	-	0.5840
Tiny ImageNet 200	MCR^2	227.6468	0.1319
	V- MCR^2	231.1538	0.2665
	CE	-	0.1907

Inner product of representations (Left)

Comparison of classification performance (Center)

Note (left) : Sort the columns of Z_θ by classes and calculate absolute value of inner product $|Z_\theta^T Z_\theta|$

$V - MCR^2$ related thoughts

- Although $V - MCR^2$ is approximation of MCR^2 , it shows better representation learning performance in the criteria : (seemingly more direct criteria compared to InfoMax principle)

high quality representation \Leftrightarrow if points from different classes lie on separate, orthogonal subspaces, and the union these subspaces span as many dimensions as possible.

1. What optimization process result in better representation learning ?
 - The common phenomenon on $V - MCR^2$ is that R^c term is further decreased compared to MCR^2 , which directly contributed the orthogonality of subspaces
 - It seems that reparameterization plays key role in further minimizing the R^c term.

$V - MCR^2$ related thoughts

2. The original MCR^2 is heading for low dimensional **linear** subspaces that classifies well.
 - Hence, it will obviously improve downstream performance under linear evaluation protocol if the learned subspaces reflect the representation well enough.

3. Again, **the assumption that $z \sim MN$ or $Z \sim \text{mixture of } MN$ is not guaranteed here.**
 - If the assumption holds, it would be better to just use EM algorithm on learned representation. (EM algorithm is well-known for good performance on grouping mixture of MN)
 - To avoid this assumption, we may be able to use variational bounds for entropy of $p(z)$ to replace R or R_c term. ($R \sim H(Z)$, $R_c \sim H(Z|Y)$)
 - Before this, why this assumption can be assumed in very complex data set???
 - Or we can design our architecture for f_θ to enforce the learned representation to follow mixture of MN.