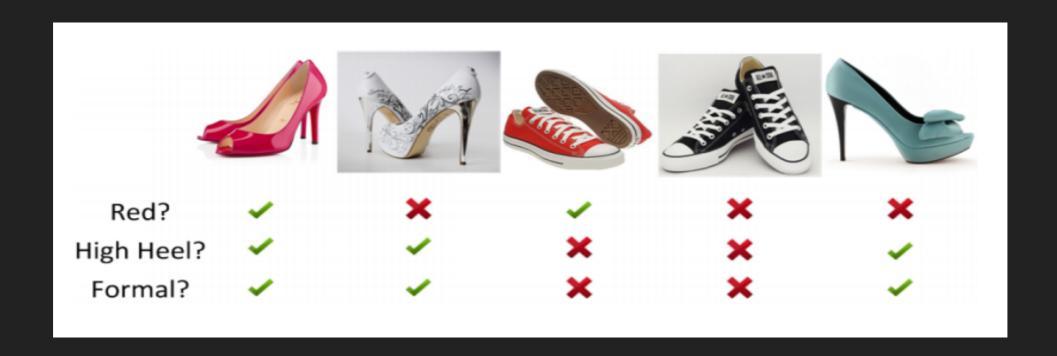
# Clustering-Based Joint Feature Selection for semantic attribute prediction



#### PRESENTED BY: TEAM GENESIS

[AVANI RAWAT] - 20172029

[SURBHI GOYAL] - 20172023

[SAMEER SINGH] - 20172134

[SHUBHAM VERMA] - 20172035

# **MOTIVATION**

- Need to handle High Dimensional Data
- Representation of semantic features effectively
- Effective utilisation of relatedness among multiple attributes
- Not all low-level features have equal contribution to all the attributes
- Attributes are related in clustering structures

#### INTRODUCTION

- Features: Metric for property or characteristic of an entity
- Semantic Features: Give more understanding about the entity that can be used for its recognition/classification
- AIM: To learn semantic features
- However, for this learning if we use raw high-dimensional data, it will suffer from the curse of dimensionality
- Besides, not all low-level features have equal contribution to all the attributes

# **KEY IDEA**

- Same types of entities share some common set of features
- ▶ Thus, attributes occur in clustering structures
- Methodology:

- 
$$F = \{f1, f2, ..., fd\}$$

- Features

$$-X = \{x1, x2, ..., xn\}$$

- Data

$$- C = \{c1, c2, ..., cm\}$$

- Classes

$$-Y = {y1,y2,...,yn}$$

- Predicted Labels

$$-s = f(0, ..., 0, 1, ..., 1)$$

- Feature Selection Vector

$$-W = [w1, w2, ..., wm]$$

- Linear Projection Matrix for mapping X to Y

$$\min_{W,\mathbf{s}} \ L(W^{\top} \operatorname{diag}(\mathbf{s}) X, Y)$$
 
$$s.t., \ \mathbf{s} \in \{0,1\}^n, \ \mathbf{s}^T \mathbf{1}_n = K$$

# FEATURE SELECTION VECTOR

$$s = f(0,...,0,1,...,1)$$

- 'f' is a permutation function that contains 0 and 1s
- S is the feature selection vector that is a binary vector
- It contains 1 for the kth cluster as selected
- And, 0 for the kth clustered not selected

# LINEAR PROJECTION MATRIX

$$W_i = [\boldsymbol{w}_1^{(i)}, \boldsymbol{w}_2^{(i)}, \dots, \boldsymbol{w}_{n_i}^{(i)}]$$

- Linear Projection Matrix contains the Linear Projection Vectors Wi
- Where, Wi = [wi(1),wi(2),...,wi(ni)]
- where, wi(1) is the weight vector for the ith cluster

#### STEP TOWARDS LABEL CORRELATION

- The linear projection matrix that was described, doesn't consider any correlation among their vectors
- Correlation is captured using the mean vector for ith cluster and the difference is tracked with each weight vector
- K-means is utilized to capture the correlation

# MODELING LABEL CORRELATION

- Idea: Correlated attributes share the same features
- Model: Learning the clustering structures through k-means
- ▶ WE = [W1,W2,...,Wk]

where, 
$$\mathbf{W} = [w_1^{(i)}, w_2^{(i)}, \dots, w_{n_i}^{(i)}]$$

and,  $\mathbf{E} = \text{Permutation Partition Matrix}^{\dagger}$ 

$$\sum_{i=1}^{k} \sum_{j=1}^{n_i} \|\boldsymbol{w}_j^{(i)} - \boldsymbol{m}_i\|^2, \boldsymbol{m}_i = \sum_{j=1}^{n_i} \boldsymbol{w}_j^{(i)} / n_i$$
 (1)

#### DEDUCTION OF COST FUNCTION IN VECTORIZED FORM

$$\sum_{i=1}^{k} \sum_{j=1}^{n_i} \|\boldsymbol{w}_j^{(i)} - \boldsymbol{m}_i\|^2 = \sum_{i=1}^{k} \|W_i(I_{n_i} - \frac{e_i e_i^{\top}}{n_i})\|_F^2$$

$$= \sum_{i=1}^{k} \mathbf{Tr}(W_i^{\top} W_i) - (\frac{e_i^{\top}}{\sqrt{n_i}}) W_i^{\top} W_i (\frac{e_i}{\sqrt{n_i}})$$
 (2)

Let  $F = \text{diag}(\frac{e_1}{\sqrt{n_1}}, \frac{e_2}{\sqrt{n_2}}, \dots, \frac{e_k}{\sqrt{n_k}}) \in \mathbb{R}^{m \times k}$  be an orthonormal matrix, then Eq. (2) can be rewritten as

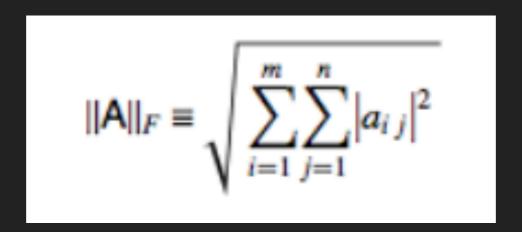
$$\mathbf{Tr}(W^{\top}W) - \mathbf{Tr}(F^{\top}W^{\top}WF)$$

mi = Mean Vector of ith-cluster

$$ei = transpose([1,1,...,1])$$
 (ni x 1) vector

# FROBENIUS NORM

- The Frobenius norm is also called the Euclidean norm
- Frobenius norm is matrix norm of an m×n matrix A defined as the square root of the sum of the absolute squares of its elements



# REDUCED OPTIMIZATION PROBLEM

$$\min_{F^{\top}F=I_k} \mathbf{Tr}(W^{\top}W) - \mathbf{Tr}(F^{\top}W^{\top}WF) + \gamma \mathbf{Tr}(W^{\top}W)$$
 (3)

- Given above the optimization problem, that captures the label correlation
- The above problem needs to be solved with the feature selection model to give the overall functionality for predicting the semantic features out of the dataset

#### FEATURE SELECTION

$$\min_{W,F,\mathbf{s}} L(W^{\top} \operatorname{diag}(\mathbf{s})X, Y) + \gamma \mathbf{Tr}(W^{\top}W) 
+ \beta (\mathbf{Tr}(W^{\top}W) - \mathbf{Tr}(F^{\top}W^{\top}WF)) 
s.t. F^{\top}F = I_k, \mathbf{s} \in \{0, 1\}^n, \mathbf{s}^{\top}\mathbf{1}_n = K$$
 (4)

$$\min_{W,F;F^{\top}F=I_{k}} L(W^{\top}X,Y) + \alpha \sum_{i=1}^{k} ||W_{i}||_{2,1} + \gamma \mathbf{Tr}(W^{\top}W) + \beta (\mathbf{Tr}(W^{\top}W) - \mathbf{Tr}(F^{\top}W^{\top}WF)) \tag{5}$$

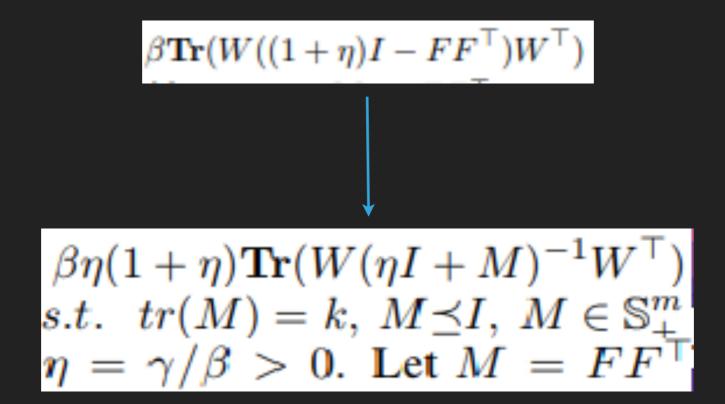
 As we have captured the attribute correlation from eqn 3, the above optimization problem combines it with the methodology

where,

- $\alpha$  = parameter to control the sparsity of W
- $\beta$  = control contribution from label correlation
- $\gamma$  = control generalization performance

#### **OPTIMIZATION**

- The optimization problem that was just deducted is non-convex, non-smooth
- Thus, the problem becomes untractable
- Relaxations need to be done to transfer it to convex one



#### ASO (ALTERNATING STRUCTURE OPTIMIZATION)

- Multi-task learning (MTL) learns multiple related tasks simultaneously to improve generalization performance
- It's the approach that is based on some common structure among various different problems
- For every problem there are 2 parts:
  - Common Property
  - Associated Property

# OPTIMIZING M WHEN FIXING W

To optimize M and W simultaneously is difficult, however using ASO it can be done easily as follows:

$$\min_{M} \mathbf{Tr}(W(\eta I + M)^{-1}W^{\top})$$
s.t.  $tr(M) = k, \ M \preceq I, \ M \in \mathbb{S}_{+}^{m}$ 

Wis decomposed using SVD while M using Eigen decomposition, and putting  $Q^* = V$ 

, 2011]

$$\Lambda^* = \arg\min_{\Lambda} \sum_{i=1}^{q} \frac{\sigma_i^2}{\eta + \lambda_i}$$

$$s.t. \quad \sum_{i=1}^{q} \lambda_i = k, 0 \le \lambda_i \le 1$$

#### OPTIMIZING W WHEN FIXING M

- L 2,1 norm is difficult to calculate
- Using dummy variables, L 2,1 is opted out by introducing dummy variables

$$\sum_{i=1}^k (\|W_i\|_{2,1})^2 = (\sum_{i=1}^k \sum_{j=1}^d \|\boldsymbol{w}_{i,j}\|_2)^2 \le \sum_{i=1}^k \sum_{i=1}^d \frac{(\|\boldsymbol{w}_{i,j}\|_2)^2}{\delta_{ij}}$$

$$\sum_{i}\sum_{j}\delta_{ij}=\overline{1}$$
  
 $\delta_{ij}\in\mathbb{R}^+$ 

where  $w_{i,j} \in \mathbb{R}^{1 \times m}$  is the row vector of  $W_i$ . Thus  $\delta_{ij}$  can be updated by holding the equality:

Where,

$$\delta_{ij} = \|\mathbf{w}_{i,j}\|_2 / \sum_{j=1}^d \|\mathbf{w}_{i,j}\|_2.$$

$$\delta_{ij} = \|\boldsymbol{w}_{i,j}\|_2 / \sum_{j=1}^{d} \|\boldsymbol{w}_{i,j}\|_2. \quad \arg\min_{W} \|W^T X - Y\|_F^2 + \alpha \sum_{i=1}^{k} \sum_{i=1}^{d} \frac{(\|\boldsymbol{w}_{i,j}\|_2)^2}{\delta_{ij}} - \beta \eta (1 + \eta) \mathbf{Tr}(W(\eta I + M)^{-1} W^\top)$$

#### FEATURE SELECTION OPTIMIZATION

#### Algorithm 1 Feature Selection Optimization

#### Input:

- Multiple attribute data {X, Y};
- Parameters α, β, k(optional) and the number of selected features K;
- The initial projection matrix W<sub>0</sub>;

#### Procedure:

- Set W = W<sub>0</sub>;
- 2: repeat
- Update M according to Eq. (8);
- Update r according to Alg. 2;
- Update δ according to Eq. (10);
- Update W according to Eq. (11);
- 7: until Converges
- Sort each feature according to ||w<sup>i</sup>||<sub>2</sub> in descending order of each group;
- return The group-wise top-K ranked features;

#### CLUSTER ASSIGNMENT ESTIMATION

# Algorithm 2 Cluster Assignment Estimation

#### Input: M;

#### Procedure:

- Approximate F by top-ranked eigenvector of Q;
- Calculate R<sub>11</sub>, R<sub>12</sub> by applying QR decomposition with column pivoting on F by Eq. (12);
- Calculate R by Eq. (13);
- calculate r by Eq. (14) for each attribute;
- return Cluster assignment vector r;

#### **ESTIMATING ATTRIBUTE ASSIGNMENT**

- The group-wise feature selection is conducted by the clustering structure (F) of the attribute
- For reconstructing F, eigen decomposition is being used

$$F^T = \underbrace{[t_{11}v_1, \cdots, t_{1s_1}v_1, \cdots, t_{k1}v_k, \cdots, t_{ks_1}v_k]}_{cluster1}$$
 where  $V^T = [v_1, v_2, \cdots, v_k] \in \mathbb{R}^{k \times k}$  is an orthogonal matrix.

# KEY IDEA FOR FEATURE SELECTION

- Using the clustering regularizer to partition the tasks into groups where strong correlation exists among tasks in the same group
- So, feature selection based on such group structures would make sure appropriate feature subsets are selected to represent the respective semantic attributes



# Thank You