

# Chapter 7

## Learning Discriminative Models in General

supplementary slides to  
*Machine Learning Fundamentals*  
© **Hui Jiang 2020**  
published by Cambridge University Press

August 2020



# Outline

- 1 A General Framework to Learn Discriminative Models
- 2 Ridge Regression and LASSO
- 3 Matrix Factorization
- 4 Dictionary Learning

# Revisit Soft SVMs

- review soft SVM formulation:

## SVM3

$$\min_{\mathbf{w}, b, \xi_i} \quad \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_{i=1}^N \xi_i$$

subject to

$$y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i \quad \forall i \in \{1, 2, \dots, N\}$$

$$\xi_i \geq 0 \quad \forall i \in \{1, 2, \dots, N\}$$

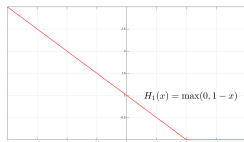
- reformulate the objective function:

$$\xi_i^* = H_1(y_i(\mathbf{w}^\top \mathbf{x}_i + b)) \quad \forall i \in \{1, 2, \dots, N\}$$

$$\begin{cases} \xi_i \geq 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b) \\ \xi_i \geq 0 \end{cases}$$

$$\Rightarrow \xi_i \geq \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b))$$

introduce the hinge function  $H_1(\cdot)$ :



$$\xi_i \geq H_1(y_i(\mathbf{w}^\top \mathbf{x}_i + b))$$

# Learning Discriminative Models

- soft SVMs can be reformulated as:

$$\min_{\mathbf{w}, b} \left[ \underbrace{\sum_{i=1}^N H_1(y_i(\mathbf{w}^\top \mathbf{x}_i + b))}_{\text{empirical loss}} + \underbrace{\lambda \cdot \|\mathbf{w}\|^2}_{\text{regularization term}} \right]$$

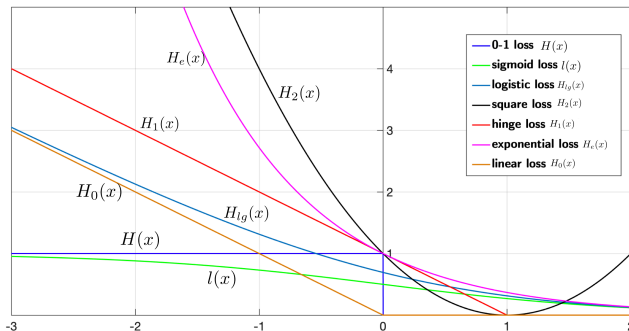
- the empirical loss is summed over all training samples when evaluated using the **hinge loss** function
  - the **regularization term** is the  $L_2$  norm of model parameters
- a general way to learn discriminative models is to minimize:

empirical loss + regularization term

# Loss Functions in Machine Learning (1)

ML method	loss function
-	0-1 loss: $H(x) = \begin{cases} 1 & x \leq 0 \\ 0 & x > 0 \end{cases}$
<b>Perceptron</b>	rectified linear loss: $H_0(x) = \max(0, -x)$
<b>MCE</b>	sigmoid loss: $l(x) = \frac{1}{1+e^x}$
<b>Logistic Regression</b>	logistic loss: $H_{\text{lg}}(x) = \ln(1 + e^{-x})$
<b>Linear Regression</b>	square loss: $H_2(x) = (1 - x)^2$
<b>Soft SVM</b>	hinge loss: $H_1(x) = \max(0, 1 - x)$
<b>Boosting</b>	exponential loss: $H_e(x) = e^{-x}$

# Loss Functions in Machine Learning (2)



- 1 monotone non-increasing
  - ex. *square loss*
- 2 convex
  - ex. *0-1 loss*
  - ex. *sigmoid loss*
- 3 sharply increasing as  $x \rightarrow -\infty$ 
  - e.g. *square loss*
  - e.g. *exponential loss*

# Regularization in Machine Learning

- soft SVMs can be formulated as the constrained optimization

$$\min_{\mathbf{w}, b} \sum_{i=1}^N H_1(y_i(\mathbf{w}^\top \mathbf{x}_i + b))$$

subject to

$$\|\mathbf{w}\|^2 \leq 1$$

- regularization  $\implies$  constraining model space in learning
- $L_p$  norm regularization ( $\forall p \geq 0$ )

$$\|\mathbf{w}\|_p \leq 1$$

where

$$\|\mathbf{w}\|_p = \left( |w_1|^p + |w_2|^p + \cdots + |w_n|^p \right)^{\frac{1}{p}}$$

# $L_p$ norm (1)

- $L_2$  norm:

$$\|\mathbf{w}\|_2 = \sqrt{|w_1|^2 + \cdots + |w_n|^2}$$

- $L_1$  norm:

$$\|\mathbf{w}\|_1 = |w_1| + \cdots + |w_n|$$

- $L_0$  norm:

$$\|\mathbf{w}\|_0 = |w_1|^0 + \cdots + |w_n|^0$$

- $\|\mathbf{w}\|_0$  ( $\in \mathbb{Z}$ ) equals the number of non-zero elements in  $\mathbf{w}$

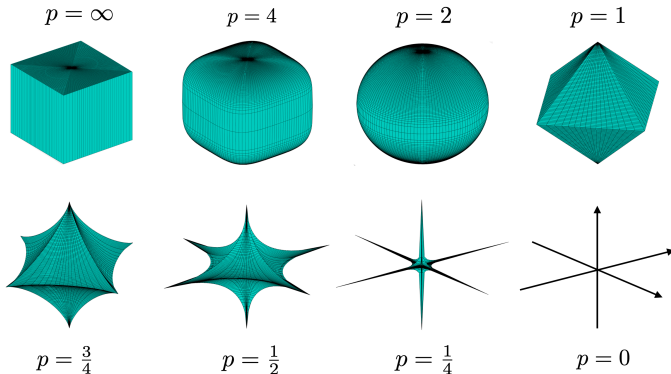
- $L_\infty$  norm:

$$\|\mathbf{w}\|_\infty = \max(|w_1|, \cdots, |w_n|)$$

- $\|\mathbf{w}\|_\infty$  equals to the largest magnitude of all elements in  $\mathbf{w}$



# $L_p$ norm (2)



- $p \downarrow \implies$   
stronger regularization
- $p \geq 1 \iff$   
convex set
- $p = 1$  is  
important

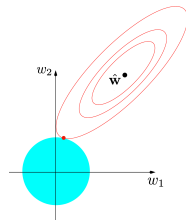
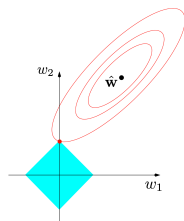
# $L_1$ Norm Regularization Promotes Sparsity

- $L_1$  norm leads to sparse solutions
- the gradient of  $L_1$  norm is constant until the parameter becomes zero

$$\frac{\partial \|\mathbf{w}\|_1}{\partial w_i} = \text{sgn}(w_i) = \begin{cases} 1 & w_i > 0 \\ 0 & w_i = 0 \\ -1 & w_i < 0 \end{cases}$$

- the gradient of  $L_2$  norm shrinks as the parameter becomes small

$$\frac{\partial \|\mathbf{w}\|_2^2}{\partial w_i} = 2w_i$$



# Ridge Regression

- ridge regression = linear regression +  $L_2$  norm regularization
- given a training set as  $\mathcal{D} = \{(\mathbf{x}_i, y_i) \mid i = 1, 2, \dots, N\}$

$$\mathbf{w}_{\text{ridge}}^* = \arg \min_{\mathbf{w}} \left[ \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \lambda \cdot \|\mathbf{w}\|_2^2 \right]$$

- the closed form solution:

$$\mathbf{w}_{\text{ridge}}^* = \left( \mathbf{X}^\top \mathbf{X} + \lambda \cdot \mathbf{I} \right)^{-1} \mathbf{X}^\top \mathbf{y}$$

- may use gradient descent to avoid the matrix inversion

# LASSO

- LASSO = linear regression +  $L_1$  norm regularization
- given a training set as  $\mathcal{D} = \{(\mathbf{x}_i, y_i) \mid i = 1, 2, \dots, N\}$

$$\mathbf{w}_{\text{lasso}}^* = \arg \min_{\mathbf{w}} \underbrace{\left[ \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \lambda \cdot \|\mathbf{w}\|_1 \right]}_{Q_{\text{lasso}}(\mathbf{w})}$$

- no closed-form solution exists, need to use gradient descent:

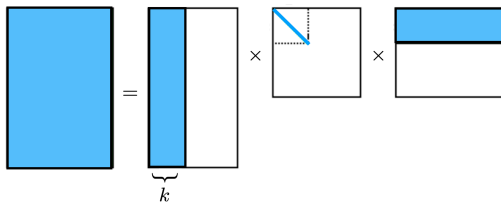
$$\frac{\partial Q_{\text{lasso}}(\mathbf{w})}{\partial \mathbf{w}} = \left( \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^\top \right) \mathbf{w} - \sum_{i=1}^N y_i \mathbf{x}_i + \lambda \cdot \text{sgn}(\mathbf{w})$$

- LASSO imposes stronger  $L_1$  regularization and gives sparse solutions, suitable for feature selection and interpretation

# Matrix Factorization (MF): SVD

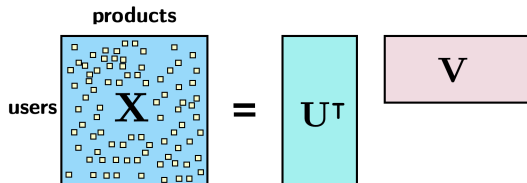
- use singular value decomposition (SVD) to factorize a matrix
- truncate to approximate:  $[\mathbf{X}]_{n \times m} \approx [\mathbf{U}]_{n \times k} [\mathbf{\Sigma}]_{k \times k} [\mathbf{V}]_{k \times m}$
- matrix factorization leads to many real-world applications

$$[\mathbf{X}]_{n \times m} = [\mathbf{U}]_{n \times m} [\mathbf{\Sigma}]_{m \times m} [\mathbf{V}]_{m \times m}$$



$$[\mathbf{X}]_{n \times m} \approx [\mathbf{U}]_{n \times k} [\mathbf{\Sigma}]_{k \times k} [\mathbf{V}]_{k \times m}$$

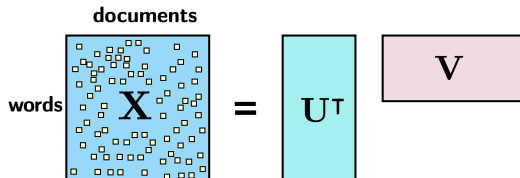
# MF Application (1): Collaborative Filtering



The diagram shows a large blue square matrix labeled  $X$  with 'users' on the left and 'products' on top. The matrix is sparse, with small yellow squares representing non-zero entries. To the right of  $X$  is an equals sign, followed by a light blue vertical rectangle labeled  $U^T$ , and then a light pink horizontal rectangle labeled  $V$ .

- collaborative filtering is the key technique for recommendation
- rely on factorizing a sparse matrix into two small dense matrices
  - construct the so-call user-product matrix
  - yielding product vectors and user vectors
  - measure similarity between different products (or users)

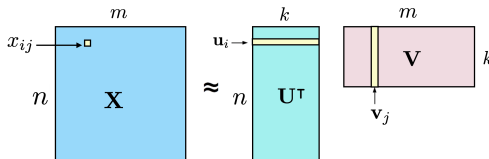
# MF Application (2): Latent Semantic Analysis



- latent semantic analysis (LSA) is an important technique in natural language processing
- rely on factorizing a sparse matrix into two small dense matrices
  - construct the so-call document-word matrix
  - yielding word vectors and document vectors
  - measure similarity between different words (or documents)

# Matrix Factorization as Machine Learning

- SVD is not efficient for large sparse matrices, not suitable for partially observed matrices
- cast matrix factorization as a machine learning problem



- learn  $\mathbf{U}$  and  $\mathbf{V}$  to minimize the reconstruction error for all observed elements in  $\Omega$  and some regularization terms:

$$Q(\mathbf{U}, \mathbf{V}) = \sum_{(i,j) \in \Omega} (x_{ij} - \mathbf{u}_i^T \mathbf{v}_j)^2 + \lambda_1 \sum_{i=1}^n \|\mathbf{u}_i\|_2^2 + \lambda_2 \sum_{j=1}^m \|\mathbf{v}_j\|_2^2$$



# Alternating Algorithm for Matrix Factorization

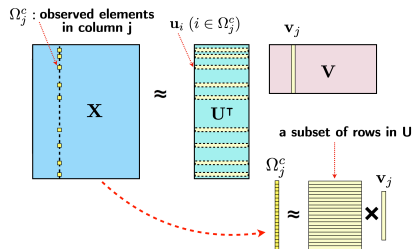
bilinear models suggest an alternating algorithm

1. keep  $\mathbf{U}$  constant,  
estimate all  $\mathbf{v}_j$  in  $\mathbf{V}$

$$\mathbf{v}_j = \left( \sum_{i \in \Omega_j^c} \mathbf{u}_i \mathbf{u}_i^\top + \lambda_2 \mathbf{I} \right)^{-1} \left( \sum_{i \in \Omega_j^c} x_{ij} \mathbf{u}_i \right)$$

2. keep  $\mathbf{V}$  constant,  
estimate all  $\mathbf{u}_j$  in  $\mathbf{U}$

$$\mathbf{u}_i = \left( \sum_{j \in \Omega_i^r} \mathbf{v}_j \mathbf{v}_j^\top + \lambda_1 \mathbf{I} \right)^{-1} \left( \sum_{j \in \Omega_i^r} x_{ij} \mathbf{v}_j \right)$$



$$\arg \min_{\mathbf{v}_j} \sum_{i \in \Omega_j^c} (x_{ij} - \mathbf{u}_i^\top \mathbf{v}_j)^2 + \lambda_2 \cdot \|\mathbf{v}_j\|_2^2$$

# Alternating Algorithm for Matrix Factorization

set  $t = 0$

randomly initialize  $\mathbf{v}_j^{(0)}$  ( $j = 1, 2, \dots, m$ )

**while** not converged **do**

**for**  $i = 1, \dots, n$  **do**

$$\mathbf{u}_i^{(t+1)} = \left( \sum_{j \in \Omega_i^r} \mathbf{v}_j^{(t)} (\mathbf{v}_j^{(t)})^\top + \lambda_1 \mathbf{I} \right)^{-1} \left( \sum_{j \in \Omega_i^r} x_{ij} \mathbf{v}_j^{(t)} \right)$$

**end for**

**for**  $j = 1, \dots, m$  **do**

$$\mathbf{v}_j^{(t+1)} = \left( \sum_{i \in \Omega_j^c} \mathbf{u}_i^{(t+1)} (\mathbf{u}_i^{(t+1)})^\top + \lambda_2 \mathbf{I} \right)^{-1} \left( \sum_{i \in \Omega_j^c} x_{ij} \mathbf{u}_i^{(t+1)} \right)$$

**end for**

$t = t + 1$

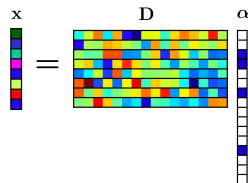
**end while**

- all updates in each step are parallelizable
- it may be faster to use SGD for very sparse matrices

# Dictionary Learning (I)

- **sparse coding assumption:** each real sample is constructed from a large dictionary based on a sparse code

$$\mathbf{x} = \begin{bmatrix} | & & | \\ \mathbf{d}_1 & \cdots & \mathbf{d}_n \\ | & & | \end{bmatrix}_{d \times n} \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix} = \mathbf{D} \boldsymbol{\alpha}$$



- **dictionary learning:** jointly learn the dictionary  $\mathbf{D}$  and sparse codes from  $\{\boldsymbol{\alpha}_i\}$  from some training samples  $\{\mathbf{x}_i\}$

$$\mathbf{X} = \begin{bmatrix} | & & | \\ \mathbf{x}_1 & \cdots & \mathbf{x}_N \\ | & & | \end{bmatrix}_{d \times N} \quad \mathbf{A} = \begin{bmatrix} | & & | \\ \boldsymbol{\alpha}_1 & \cdots & \boldsymbol{\alpha}_N \\ | & & | \end{bmatrix}_{n \times N}$$

# Dictionary Learning (II)

- dictionary learning is cast as a machine learning problem:
  - minimize the reconstruction error
  - add  $L_1$  regularization to impose sparsity of all codes
  - add  $L_2$  regularization on dictionary to avoid overfitting

$$\arg \min_{\mathbf{D}, \mathbf{A}} \underbrace{\frac{1}{2} \sum_{i=1}^N \left\| \mathbf{x}_i - \mathbf{D} \boldsymbol{\alpha}_i \right\|_2^2 + \lambda_1 \sum_{i=1}^N \left\| \boldsymbol{\alpha}_i \right\|_1 + \frac{\lambda_2}{2} \sum_{j=1}^n \left\| \mathbf{d}_j \right\|_2^2}_{Q(\mathbf{D}, \mathbf{A})}$$

- compute all gradients:

$$\frac{\partial Q(\mathbf{D}, \mathbf{A})}{\partial \mathbf{A}} = \mathbf{D}^\top \mathbf{D} \mathbf{A} - \mathbf{D}^\top \mathbf{X} + \lambda_1 \cdot \text{sgn}(\mathbf{A})$$

$$\frac{\partial Q(\mathbf{D}, \mathbf{A})}{\partial \mathbf{D}} = \mathbf{D} \mathbf{A} \mathbf{A}^\top - \mathbf{X} \mathbf{A}^\top + \lambda_2 \cdot \mathbf{D}$$

# Dictionary Learning (III)

## Gradient Descent for Dictionary Learning

set  $t = 0$  and  $\eta_0$

randomly initialize  $\mathbf{D}^{(0)}$  and  $\mathbf{A}^{(0)}$

**while** not converged **do**

    update  $\mathbf{A}$ :

$$\mathbf{A}^{(t+1)} = \mathbf{A}^{(t)} - \eta_t \left( (\mathbf{D}^{(t)})^\top \mathbf{D}^{(t)} \mathbf{A}^{(t)} - (\mathbf{D}^{(t)})^\top \mathbf{X} + \lambda_1 \cdot \text{sgn}(\mathbf{A}^{(t)}) \right)$$

    update  $\mathbf{D}$ :

$$\mathbf{D}^{(t+1)} = \mathbf{D}^{(t)} - \eta_t \left( \mathbf{D}^{(t)} \mathbf{A}^{(t+1)} (\mathbf{A}^{(t+1)})^\top - \mathbf{X} (\mathbf{A}^{(t+1)})^\top + \lambda_2 \cdot \mathbf{D}^{(t)} \right)$$

    adjust  $\eta_t \rightarrow \eta_{t+1}$

$t = t + 1$

**end while**

# Sparse Coding

- sparse coding: given a dictionary  $\mathbf{D}$ , find the sparse code  $\alpha$  for a new observation  $\mathbf{x}$
- ideal but infeasible:

$$\arg \min_{\alpha} \|\alpha\|_0 \quad \text{s. t.} \quad \mathbf{D} \alpha = \mathbf{x}$$

- a practical solution:
  - replace intractable  $L_0$  norm with  $L_1$
  - relax to an imperfect reconstruction

$$\alpha^* = \arg \min_{\alpha} \underbrace{\frac{1}{2} \|\mathbf{x} - \mathbf{D} \alpha\|_2^2 + \lambda_1 \cdot \|\alpha\|_1}_{Q'(\alpha)}$$

- use gradient descent:

$$\frac{\partial Q'(\alpha)}{\partial \alpha} = \mathbf{D}^T \mathbf{D} \alpha - \mathbf{D}^T \mathbf{x} + \lambda_1 \cdot \text{sgn}(\alpha)$$