

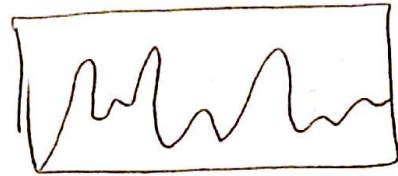
Supervised Dictionary Learning

report By :

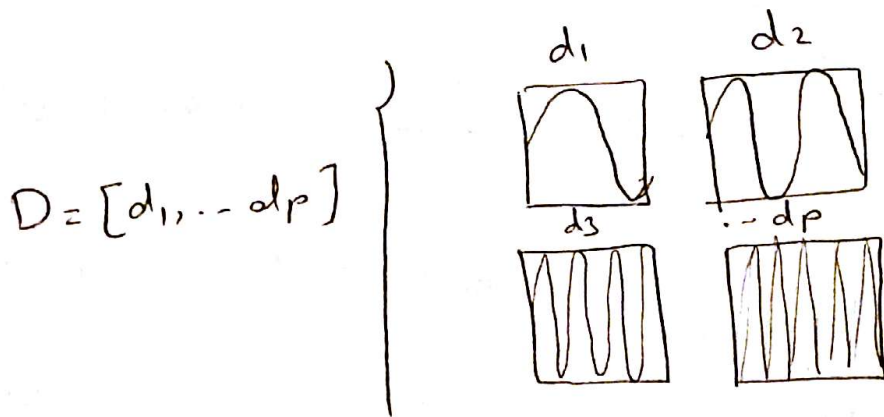
Mohsen Nabian

$$\underbrace{y}_{\text{measured image}} = \underbrace{X_{\text{orig}}}_{\text{original image}} + \underbrace{\omega}_{\text{noise}}$$

Let x in \mathbb{R}^m be a signal:



Let $D = [d_1, \dots, d_p] \in \mathbb{R}^{m \times p}$ be a set of normalized "basis" vectors we call it dictionary.



D is adapted to y if it can represent it with a few basis vector - that is, there exist a sparse vector α , in \mathbb{R}^p such that $y \approx D\alpha$.
we call α the sparse code.

$\alpha \in \mathbb{R}^p$
sparse

$$\begin{pmatrix} y \end{pmatrix} \approx \begin{pmatrix} d_1 & d_2 & \dots & d_p \end{pmatrix} \begin{pmatrix} \alpha[1] \\ \alpha[2] \\ \vdots \\ \alpha[p] \end{pmatrix}$$

4

let $(\hat{y}^i, \hat{x}^i)_{i=1}^n$ be a training set, where vector \hat{x}^i are in \mathbb{R}^p and called features. The scalars \hat{y}^i are in

$\hat{y}^i \in \begin{cases} \{-1, +1\} & \text{for binary class classification} \\ \{1, \dots, N\} & \text{for multiclass} \\ \mathbb{R} & \text{for regression problems} \end{cases}$

in linear model: $y \approx w^T x$ or $y \approx \text{sign}(w^T x)$

and solve: $\min_{w \in \mathbb{R}^p} \underbrace{\frac{1}{n} \sum_{i=1}^n \ell(\hat{y}^i, w^T \hat{x}^i)}_{\text{data fitting}} + \underbrace{\lambda \Omega(w)}_{\text{regularization}}$

Optimization for Dictionary learning

$$\min_{\substack{\alpha \in \mathbb{R}^{p \times n} \\ D \in C}} \left[\frac{1}{2} \|y_i - D \alpha_i\|_2^2 + \lambda \|\alpha_i\|_1 \right] \quad (*)$$

$$C \triangleq \{ D \in \mathbb{R}^{m \times p} \text{ s.t. } \forall j=1, \dots, p, \|d_{\cdot j}\|_2 \leq 1 \}$$

•) classical optz. alternate between D & α . but they are slow.

- we can represent each data as a ~~Comb~~ sparse Linear Combination of dictionary vectors: assume we know dictionary D and x given \Rightarrow we have to find sparse Coefficients α :

$$R^*(x, D)$$

$$R^*(x, D) = \min_{\alpha \in \mathbb{R}^k} \|x - D\alpha\|_2^2 + \lambda_1 \|\alpha\|_1$$

in this paper, we consider:

- each signal x_i belongs to a class $y_i \in \{y_1, y_2, \dots, y_p\}$
- Consider $p=2 \Rightarrow y_i \in \{-1, +1\}$
- goal: learn jointly a single dictionary D adapted to classification task and function f which should be positive for each signal in class $+1$ and negative otherwise.

model i) [linear in α]: $f(x, \alpha, \theta) = w^T \alpha + b$
 (L) where $\theta = \{w \in \mathbb{R}^k, b \in \mathbb{R}\}$

model ii) [bilinear in x and α]

(NL) $f(x, \alpha, \theta) = x^T W \alpha + b$
 where $\theta = \{W \in \mathbb{R}^{n \times k}, b \in \mathbb{R}\}$

model is has more parameter is richer model.

So we must ~~learn~~ ^{learn} α first and then ~~learn~~ ^{learn} θ parameters.

we classically obtain α as follows:

$$\min_{D, \alpha} \sum_{i=1}^m \|x_i - D\alpha_i\|_2^2 + \lambda_1 \|\alpha_i\|_1 \quad (2)$$

(2) will end up with sparse α_i and ~~the~~
Single dictionary D

Then we learn function parameters from supervised learning :

$$\min_{\theta} \sum_{i=1}^m C(y_i f(x_i, \alpha_i, \theta)) + \lambda_2 \|\theta\|_2^2 \quad (3)$$

$$C(x) = \log(1 + e^{-x})$$

However our goal is to learn D and θ jointly.

⇒ we propose:

$$\min_{D, \theta, \alpha} \left(\sum_{i=1}^m C(y_i, f(x_i, \alpha_i, \theta)) + \lambda_0 \|x_i - D\alpha_i\|_2^2 + \lambda_1 \|\alpha_i\|_1 + \lambda_2 \|\theta\|_2^2 \right) \quad (4)$$

SDL-G

training $\begin{matrix} \swarrow \lambda_0, \lambda_1, \lambda_2 \\ \searrow \end{matrix}$ $\begin{matrix} D \checkmark \\ \theta \checkmark \end{matrix}$

define: $S(\alpha, x_i, D, \theta, y_i) = C(y_i, f(x_i, \alpha_i, \theta)) + \lambda_0 \|x_i - D\alpha_i\|_2^2 + \lambda_1 \|\alpha_i\|_1$

$$S^*(x_i, D, \theta, y_i) = \min_{\alpha} S(\alpha, x_i, D, \theta, y_i)$$

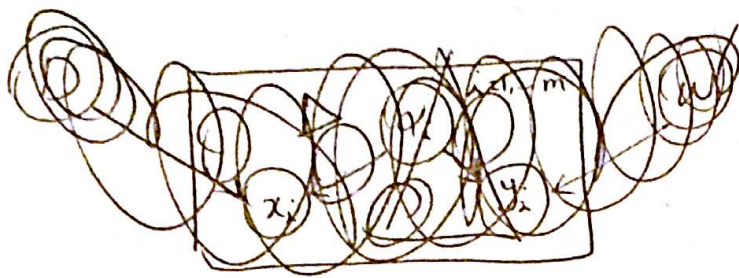
Once we obtained D, θ from (4):

given a new signal x with unknown label y

classification is as follows:

(6) $\min_{y \in \{-1, +1\}} S^*(x, D, \theta, y)$ classification of test data

Graphical Model For proposed Learning Framework:



Since ~~test~~ test classification (6) compares different Csts $S^*(x, D, \theta, y)$ we have to make sure Cst $S^*(x_i, D, \theta, -y_i)$ is larger than $S^*(x_i, D, \theta, y_i)$. Therefore we modify (4) to

$$\min \left\{ \sum_{i=1}^m \mu C(S^*(x_i, D, \theta, -y_i) - S^*(x_i, D, \theta, y_i)) + (1-\mu) S^*(x_i, D, \theta, y_i) + \lambda_2 \|\theta\|_2^2 \right\}$$

μ : trade-off coef.

(8)

SDL-D

learning jointly $\begin{matrix} D \\ \theta \end{matrix}$

using training set

we can extend (8) to multi-class classification by Soft-max.

- ① Then in the paper the algorithm for SDL-D (18) is provided which I do not write it here.
- ② ~~(6)~~ (6) which has D and θ given, is a convex optimization. The author suggest to use fixed-point-continuation-method (FPC) to solve this minimization.

Experimental Validation

Graphical Model of SDL model :

The correspondingly graphical model for our model is as follows :

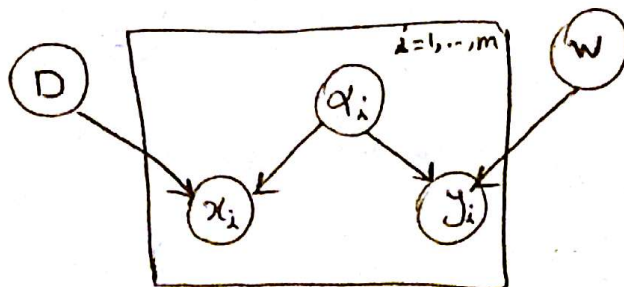


plate model

Experimental Validation :

data sets : (handwritten digits)

MNIST: 70000 28x28 image 60000 train 10000 test

USPS : 16x16 image 7291 train 2007 test

assuming $k = \frac{\lambda_1}{\lambda_0}$ K : size of dictionary

Five-fold cross validation on:

$k = \{0.13, 0.14, 0.15, 0.16, 0.17\}$

$K = \{24, 32, 48, 64, 96\}$

Classification results:

	RECL	SDL-GL	SDL-DL	RECDL	K-MN	SVM
MNIST	4.33	3.5	1.05	3.4	5.0	1.4
USPS	6.83	6.67	3.54	4.3	5.2	4.2

outperforming others.

2nd experiment:

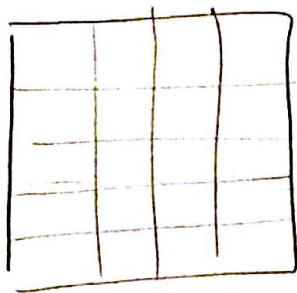
"are the obtained D discriminative per se?"

Trained USPS 10 binary classifiers (one vs all)

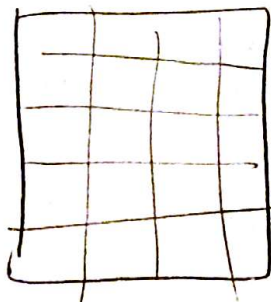
So for given value $\mu \Rightarrow$ obtain $\left\{ \begin{array}{l} 10 \text{ dictionary } D \\ \text{SDL-DL model} \\ 10 \text{ priors } \theta \end{array} \right.$

So we decompose each image \Rightarrow gives α

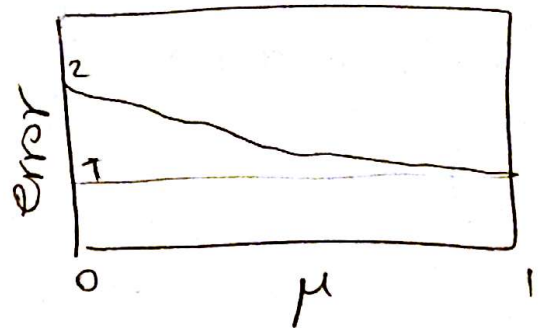
then use α in SVM. \Rightarrow error plot provided below vs $\mu \Downarrow$



(a) REC dictionaries



(b) SDL-D dictionaries



3rd experiment

Texture classification:

in experiment 1, \perp model outperformed BL model. one reason could be simplicity of

9/

the problem.

Two different textures (2 classes) are provided.
each image is broken to 12×12 patches.

all left half patches \Rightarrow train

\sim right $\sim \sim \Rightarrow$ testing

① Error rate of classification are provided.
and showed that BL models significantly
outperformed L models. and the reason is
the complexity of the problem. ✓