

Lecture 12: SSL demo

We dedicate this lecture to a demo of graphical SSL for multi-class classification of hand written digits, the MNIST dataset.

We will see:

- How to construct similarity graphs on realistic data sets.
- How to modify SSL to perform multi-class classification w/ one-hot encoding.

We will consider the problem of classifying images of handwritten digits from MNIST.

In particular we look at a subset of MNIST consisting of roughly 600 images of 1's, 4's & 8's.

But we wish to label the entire dataset w/ only 10 randomly selected labelled images!

Features	4	1	8	1	8	4	1	1
	1	4	1	1	4	4	1	8
	1	1	4	1	1	4	1	1
	1	1	4	1	4	1	8	4
	1	1	4	4	8	1	4	4
	1	4	4	8	4	8	1	4
	1	1	4	1	8	8	1	1
	4	1	1	1	8	8	4	

ie, we only have 10 labelled images & the rest are unlabelled! We will solve this prob. using the graphical SSL technique of Lec 21 that was demo'd on a toy data set.

22.1 One-hot encoding

Before getting to the algorithmic details we have to deal with the issue of multi-class classification.

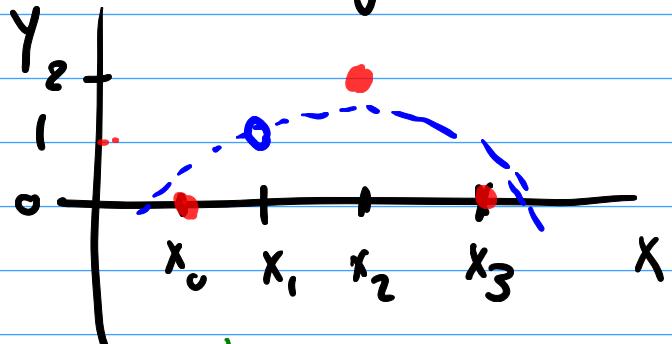
So far in the course we have mostly worked with binary classification where the two classes could be encoded as $\{-1, +1\}$.

This worked very well & was natural in many ways.

But this encoding of binary labels does not accommodate more than two classes such as our example!

A naive approach would simply be to use the encoding $\{0, 1, 2\}$ for example. Indeed, this is how the MNIST labels are given in the dataset.

But this encoding does not work well in practice & in particular when regression is employed.



essentially if you arrange the $\{0, 1, 2\}$ labels in a vector $\mathbf{y} \in \mathbb{R}^M$ & do regression on it you are implicitly introducing an ordering of the classes where your predictor cannot go straight from class $\{0\}$ to $\{2\}$ without predicting class $\{1\}$ in between!

This ordering is arbitrary & can lead to inconsistencies.

To remedy this we will use a one-hot encoding!

Rather than labeling each point $x_j \in X$ with a label $y_j \in \mathbb{R}$ we instead take $\mathbf{y}_j \in \mathbb{R}^K$ when K is # of classes in data set. More

precisely, if \underline{x}_j belongs to class $k \in \{0, \dots, K-1\}$

then $\underline{y}_j = e_k = (0, 0, \dots, \underset{k^{\text{th}} \text{ entry}}{1}, 0, 0, \dots, 0) \in \mathbb{R}^K$

In the MNIST case, since we have only 3-digits ($K=3$) then our labels have the form

$$\underline{y}_j = \begin{cases} (1, 0, 0) & \leftarrow 4's \\ (0, 1, 0) & \leftarrow 8's \\ (0, 0, 1) & \leftarrow 1's \end{cases} \quad \begin{matrix} \# \text{ of labelled} \\ \downarrow \text{pts} \\ M \times K \end{matrix} \quad \begin{matrix} \# \text{ of} \\ \text{classes} \end{matrix}$$

Thus, the output/label data $Y \in \mathbb{R}^{M \times K}$ & the class of each $\underline{x}_j \in X$ is given by the index of the largest entry in the j^{th} row of Y .

This encoding does not impose any arbitrary ordering of the classes & is invariant under reordering of the classes.

22.2 SSL on MNIST

As mentioned above, we wish to perform SSL on a subset of MNIST with $N \approx 600$ images but with only $M=10$ observed labels. There are only 3 classes in the data, the digits 1, 4, & 8.

$$\text{so, } Y = \mathbb{R}^{M \times 3} \quad \leftarrow \begin{matrix} \text{one-hot} \\ \text{encoding} \end{matrix}$$

The feature $\underline{g}_j \in \mathbb{R}^{256}$ are flattened versions of MNIST 16×16 images.

a Step 1: Construct a graph on X .

take the \underline{u}_j to be the vertices of the graph with

$$w_{ij} = \exp\left(-\frac{\|\underline{u}_i - \underline{u}_j\|^2}{2\sigma^2}\right)$$

• Step 2: Compute eigen decomposition of graph Laplacian.

$$L = D - W \quad (\text{unnormalized Laplacian})$$

& write

$$L = Q \Lambda Q^T, \quad Q = \begin{bmatrix} q_1 & q_2 & \dots & q_{N-1} \end{bmatrix}$$

• Step 3: Perform regression on labelled set.

Without loss of generality we assume

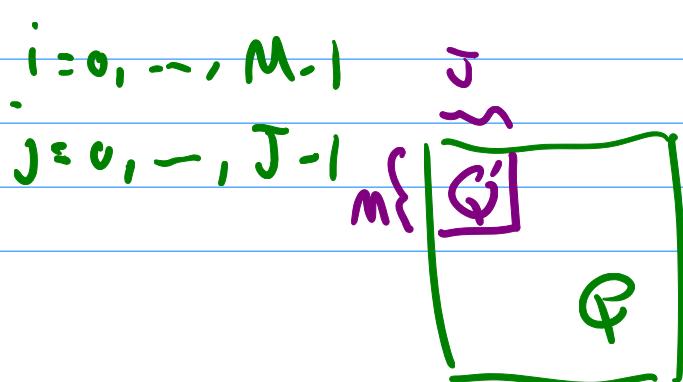
$\underline{x}_0, \dots, \underline{x}_{M-1}$ are the labelled points.

Choose $J \geq 1$ the # of eigenvectors we wish to use for regression

$Q' \in \mathbb{R}^{M \times J}$ submatrix of Q

$$Q'_{ij} = Q_{ij} \quad \text{for } i=0, \dots, M-1$$

$$j=0, \dots, J-1$$



$$\min \|Q\beta - Y\|_F^2 + \lambda \|\beta\|_F^2$$

finely solve regression problem

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{K-1} \end{bmatrix}$$

$$\hat{\beta}_k = \underset{\beta_k \in \mathbb{R}^M}{\operatorname{argmin}} \|Q\beta_k - y\|_2^2 + \lambda \|\beta_k\|_2^2$$

$$y = \begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_{K-1} \end{bmatrix}$$

* Note $\hat{\beta} \in \mathbb{R}^{J \times 3}$ since $Y \in \mathbb{R}^{M \times 3}$!

• Step 4: Predict labels on all of X.

With $\hat{\beta}$ we can predict the one-hot vectors of all pts on the graph simply as

$$\mathbb{R}^{N \times 3} \rightarrow \tilde{Y} = \begin{bmatrix} q_0 & q_1 & \dots & q_J \end{bmatrix} \hat{\beta} \in \mathbb{R}^{N \times J}$$

But we are not done yet! Rows of \tilde{Y} are not 0's & 1's. So we threshold them!

jth row
of \tilde{Y}

$$\tilde{Y}_j = \begin{bmatrix} \tilde{y}_{j0} \\ \tilde{y}_{j1} \\ \vdots \\ \tilde{y}_{jM} \end{bmatrix}$$

$$\tilde{y}_j$$

$$\tilde{y}_j \rightarrow \hat{y}_j \begin{cases} \text{set largest entry to 1} \\ \text{& rest to 0.} \end{cases}$$

- Comparing first 30 labels.
- Yellow boxes indicate the corresponding class
- error rate: 28% missclassified

