

# **MATH 156 Machine Learning**

## **Lecture Note**

Hanbaek Lyu

DEPARTMENT OF MATHEMATICS, UNIVERSITY OF CALIFORNIA, LOS ANGELES, CA 90095

*Email address:* hlyu@math.ucla.edu

WWW.HANBAEKLYU.COM

SUPPLEMENTARY REPOSITORY: [HTTPS://GITHUB.COM/HANBAEKLYU/MATH156\\_UCLA\\_SP21](https://github.com/HanbaekLyu/Math156_UCLA_SP21)

## Contents

<b>Chapter 1. Introduction to Machine Learning</b>	<b>3</b>
1. Prologue	3
2. Polynomial Regression	4
3. Decision Theory	12
<b>Chapter 2. Algorithms for Classification</b>	<b>16</b>
1. Introduction	16
2. Logistic Regression	17
3. Binary classification of MNIST by Logistic Regression	20
4. Classification of MNIST by Multiclass Logistic Regression	25
5. Probit Regression	27
6. Binary classification of MNIST by Probit Regression	29
7. Naive Bayes Classifier	31
8. Classification with the 20Newsgroups dataset	36
9. Gaussian Naive Bayes	40
<b>Chapter 3. Neural Networks</b>	<b>42</b>
1. Model formulation	42
2. Training Feedforward Neural Networks	43
3. Examples of 2-layer FFNN	49
4. Convolutional Neural Networks	51
<b>Chapter 4. Matrix Factorization and Dictionary Learning</b>	<b>58</b>
1. Introduction	58
2. Algorithms for Matrix Factorization	61
3. Applications of Matrix Factorization	62
4. Principal Component Analysis	66
5. Applications of PCA and Matrix Factorization	71
6. Probabilistic PCA and EM algorithm	73
<b>Bibliography</b>	<b>80</b>
<b>Appendix A. Review of basic probability theory</b>	<b>81</b>
1. Probability measure and probability space	81
2. Random variables	82
3. Conditional and iterated expectations	84
4. Conditional variance	85
5. Joint probability distributions	87
6. Definition and Examples of MLE	88
7. Bayesian Inference	91

## CHAPTER 1

# Introduction to Machine Learning

### 1. Prologue

When do we say we learned something, say, how to ride bikes or how to multiply two fractions? Roughly speaking, this is to *recognizing patterns* between input and output so that we can use the recognized patterns to *predict* the output of fresh new input. *Machine Learning* (ML) consists of techniques that enable computers to this learning process. Below are some key components in ML.

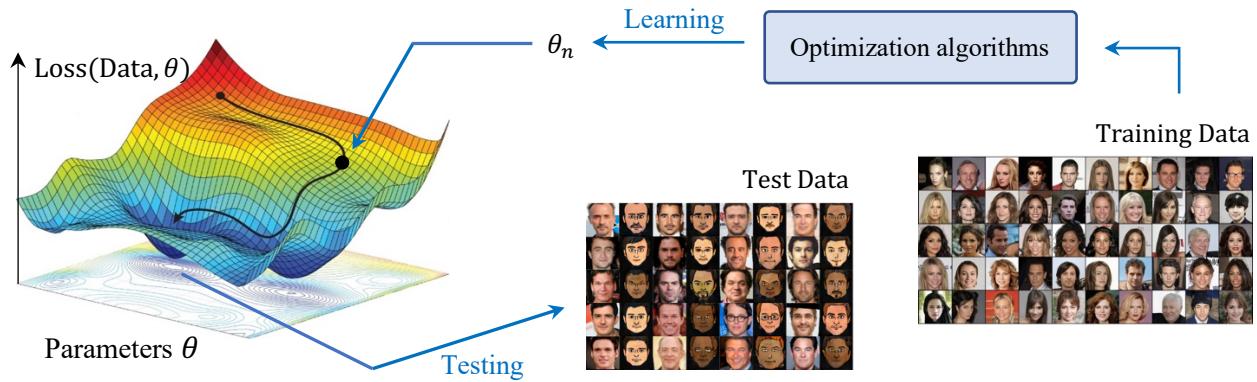


FIGURE 1. Basic scheme of Machine Learning: (1) Select a model to explain the dataset; (2) Fit the parameters to the training dataset using an optimization algorithm; (3) Test the performance of the fitted model to the test dataset.

**1. Model Selection:** The goal is to predict unknown output  $\mathbf{y}$  from a new input  $\mathbf{x}$ . For this purpose, choose a model  $\mathbf{y} \approx \hat{\mathbf{y}}(\mathbf{x}, \boldsymbol{\theta})$  with parameters  $\boldsymbol{\theta}$ <sup>1</sup>. Given input  $\mathbf{x}$  and parameters  $\boldsymbol{\theta}$ , the model will output the value  $\hat{\mathbf{y}}$  that should be close to the true and unknown  $\mathbf{y}$ .

e.g., Linear regression, Logistic regression, PCA, NMF, feedforward neural networks, CNN, RNN, etc.

**2. Training Data:** We need a large number of known input and output pairs  $(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_n, \mathbf{y}_N)$ , from which we want to recognize patterns between the variables  $\mathbf{x}$  and  $\mathbf{y}$ .

e.g., vectors, matrices, tensors, images, videos, texts, networks, etc.

**3. Training the model:** We need to find the optimal parameter  $\hat{\boldsymbol{\theta}}$  for which the model  $\mathbf{x} \mapsto \hat{\mathbf{y}}(\mathbf{x}, \hat{\boldsymbol{\theta}})$  gives a good fit for the training data. More precisely, this is done by solving following optimization problem

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmin}} \ell(\text{training data}, \boldsymbol{\theta}), \quad (1)$$

where  $\ell$  is a *loss function* that measures the error between the true output  $\mathbf{y}$  and the predicted output  $\hat{\mathbf{y}}$ , and  $\Theta$  is a set of admissible parameter values. A typical choice of such loss function is

---

<sup>1</sup>Model parameters are usually denoted by  $\theta$  in statistics and by  $w$  (for ‘weights’) in machine learning. We will use both interchangeably.

the *mean squared error* (MSE)

$$\ell(\text{training data}, \boldsymbol{\theta}) = \frac{1}{N} \sum_{k=1}^N \|\mathbf{y}_i - \hat{\mathbf{y}}(\mathbf{x}_i, \boldsymbol{\theta})\|^2, \quad (2)$$

where ‘training data’ consists of training examples  $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_N, \mathbf{y}_N)$ . Then use *optimization algorithms* to find the optimal parameter  $\hat{\boldsymbol{\theta}}$  in (1).

e.g., gradient descent, stochastic gradient descent, expectation-maximization, majorization-minimization, etc.

- 4. Testing:** After an optimal parameter  $\hat{\boldsymbol{\theta}}$  is learned, test the ‘recognized pattern’  $\mathbf{x} \mapsto \hat{\mathbf{y}}(\mathbf{x}, \hat{\boldsymbol{\theta}})$  on *test data*  $(\mathbf{x}'_1, \mathbf{y}'_1), \dots, (\mathbf{x}'_m, \mathbf{y}'_m)$ , that were kept hidden during the training step. If the loss  $\ell(\text{test data}, \hat{\boldsymbol{\theta}})$  is small, then we have learned the correct pattern successfully.

## 2. Polynomial Regression

**2.1. Basic theory.** Suppose we have a training dataset of pairs of real numbers  $(x_1, y_1), \dots, (x_N, y_N)$ . Think of  $y_k$  is the price of the stock you are holding and  $x_k$  is NASDAQ composite at time  $k$ . You want to obtain a simple functional explaining the *dependent variable*  $y_k$  in terms of the *independent variable*  $x_k$ . While there could be a general relationship between the variables  $x$  and  $y$ , here we will first study a simple model of *polynomial regression*. Namely, consider the following polynomial function  $x \mapsto \hat{y}(x; \mathbf{w})$ ,

$$\hat{y}(x; \mathbf{w}) := w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M = \sum_{k=0}^M w_k x^k = \boldsymbol{\phi}(x)^T \mathbf{w}, \quad (3)$$

where

$$\boldsymbol{\phi}(x) = [1 \ x \ \dots \ x^M]^T, \quad \mathbf{w} = [w_0 \ w_1 \ \dots \ w_M]^T \in \mathbb{R}^{M+1}. \quad (4)$$

Here  $\mathbf{w} \in \mathbb{R}^{M+1}$  denotes the  $(M+1)$ -dimensional vector of parameters. Here  $\hat{y}(x; \mathbf{w})$  is a degree  $M$  polynomial in  $x$  and the components of  $\mathbf{w}$  give the coefficients of monomials of  $x$ .

Suppose we have training data samples  $(x_1, y_1), \dots, (x_N, y_N)$ . Let  $\mathbf{Y}_{\text{train}} = [y_1, \dots, y_N]^T$ . We would like to find the optimal parameter  $\mathbf{w}$  so that the corresponding polynomial regression (3) gives a good fit for the training data. In order to do so, we first write the vector of  $n$  predicted outcomes as the following matrix form

$$\begin{bmatrix} \hat{y}(x_1, \mathbf{w}) \\ \hat{y}(x_2, \mathbf{w}) \\ \vdots \\ \hat{y}(x_N, \mathbf{w}) \end{bmatrix} = \begin{bmatrix} 1 & x_1 & \dots & x_1^M \\ 1 & x_2 & \dots & x_2^M \\ \vdots & \vdots & & \vdots \\ 1 & x_N & \dots & x_N^M \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_M \end{bmatrix} \iff \hat{\mathbf{Y}}_{\text{train}} = \mathbf{X}_{\text{train}} \mathbf{w}. \quad (5)$$

Since we want to fit the polynomial regression model in (3) to the training data, we want to choose the parameter  $\mathbf{w}$  so that

$$\mathbf{Y}_{\text{train}} \approx \hat{\mathbf{Y}} = \mathbf{X}_{\text{train}} \mathbf{w}. \quad (6)$$

In order to do so, we will use the *method of least squares*. Namely, we seek to choose  $\mathbf{w} = \hat{\mathbf{w}}$  where

$$\hat{\mathbf{w}} = \underset{\mathbf{w} \in \mathbb{R}^{M+1}}{\arg \min} \|\mathbf{Y}_{\text{train}} - \mathbf{X}_{\text{train}} \mathbf{w}\|_F^2, \quad (7)$$

where the *matrix Frobenius norm*  $\|\cdot\|_F$  is defined by  $\|A\|_F = \text{tr}(A^T A)$ . When  $A = [a_1, \dots, a_d]^T$  is a column vector, then  $A^T A$  is a scalar so that

$$\|A\|_F^2 = \text{tr}(A^T A) = A^T A = a_1^2 + a_2^2 + \dots + a_d^2, \quad (8)$$

so the Frobenius norm agrees with the Euclidean norm in this case. In Exercise (1.2.1) (with  $\mathbf{Y} = \mathbf{Y}_{\text{train}}$ ,  $\mathbf{X} = \mathbf{X}_{\text{train}}$ ,  $\mathbf{B} = \mathbf{w}$ , and  $\lambda = 0$ ), we will show that the solution of (7) is given by

$$\hat{\mathbf{w}} = (\mathbf{X}_{\text{train}}^T \mathbf{X}_{\text{train}})^{-1} \mathbf{X}_{\text{train}}^T \mathbf{Y}_{\text{train}}. \quad (9)$$

Hence for any new input  $x$ , the optimal polynomial regression model will tell us to predict the corresponding output  $y$  as

$$y \approx \hat{y}(x; \hat{\mathbf{w}}) = [1 \quad x \quad x^2 \quad \cdots \quad x^M] \hat{\mathbf{w}}. \quad (10)$$

**Exercise 1.2.1** (Method of Least Squares). Suppose we have matrices  $\mathbf{Y} \in \mathbb{R}^{d \times n}$  and  $\mathbf{X} \in \mathbb{R}^{d \times r}$ . We seek to find a matrix  $\hat{\mathbf{B}} \in \mathbb{R}^{r \times n}$  where

$$\hat{\mathbf{B}} = \underset{\mathbf{B} \in \mathbb{R}^{r \times n}}{\operatorname{argmin}} \|\mathbf{Y} - \mathbf{XB}\|_F^2 + \lambda \|\mathbf{B}\|_F^2. \quad (11)$$

Here  $\lambda \geq 0$  is called the  $L_2$ -regularization parameter. (This is an instance of unconstrained quadratic optimization problem.)

(i) Show that

$$\|\mathbf{Y} - \mathbf{XB}\|_F^2 + \lambda \|\mathbf{B}\|_F^2 = \operatorname{tr}(\mathbf{Y} - \mathbf{XB})^T (\mathbf{Y} - \mathbf{XB}) + \lambda \operatorname{tr}(\mathbf{B}^T \mathbf{B}) \quad (12)$$

$$= \operatorname{tr}(\mathbf{Y}^T \mathbf{Y}) - 2\operatorname{tr}(\mathbf{Y}^T \mathbf{XB}) + \operatorname{tr}(\mathbf{B}^T \mathbf{X}^T \mathbf{XB}) + \lambda \operatorname{tr}(\mathbf{B}^T \mathbf{B}). \quad (13)$$

(ii) Show that (use Exercise 1.2.2)

$$\frac{\partial}{\partial \mathbf{B}} \left( \|\mathbf{Y} - \mathbf{XB}\|_F^2 + \lambda \|\mathbf{B}\|_F^2 \right) = 2(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \mathbf{B} - 2\mathbf{X}^T \mathbf{Y}, \quad \frac{\partial^2}{\partial \mathbf{B}^2} \left( \|\mathbf{Y} - \mathbf{XB}\|_F^2 + \lambda \|\mathbf{B}\|_F^2 \right) = 2(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}). \quad (14)$$

(iii) From (ii), conclude that the quadratic function  $\mathbf{B} \mapsto \|\mathbf{Y} - \mathbf{XB}\|_F^2 + \lambda \|\mathbf{B}\|_F^2$  is convex, and hence its every critical point is a local minimum. (See Ref)

(iv) Suppose  $\lambda = 0$  and  $\mathbf{X}^T \mathbf{X}$  is invertible<sup>2</sup>. Then from (ii) and (iii), conclude that the quadratic function  $\mathbf{B} \mapsto \|\mathbf{Y} - \mathbf{XB}\|_F^2 + \lambda \|\mathbf{B}\|_F^2$  has a unique global minimum  $\hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ <sup>3</sup>

(v) Suppose  $\lambda > 0$ . Then argue that  $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$  is always invertible<sup>4</sup>, and the quadratic function  $\mathbf{B} \mapsto \|\mathbf{Y} - \mathbf{XB}\|_F^2 + \lambda \|\mathbf{B}\|_F^2$  has a unique global minimum  $\hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{Y}$ .

**Exercise 1.2.2** (Matrix derivatives). Show the following matrix derivatives: (Ref: [The Matrix Cookbook](#))

$$(i) \text{ (First order)} \quad \frac{\partial}{\partial \mathbf{X}} \operatorname{tr}(\mathbf{B} \mathbf{X} \mathbf{C}) = \mathbf{B}^T \mathbf{C}^T, \quad \frac{\partial}{\partial \mathbf{X}} \operatorname{tr}(\mathbf{B} \mathbf{X}^T \mathbf{C}) = \mathbf{C} \mathbf{B}.$$

$$(ii) \text{ (Second order)} \quad \frac{\partial}{\partial \mathbf{X}} \operatorname{tr}(\mathbf{X}^T \mathbf{B} \mathbf{X}) = \mathbf{B} \mathbf{X} + \mathbf{B}^T \mathbf{X}, \quad \frac{\partial}{\partial \mathbf{X}} \operatorname{tr}(\mathbf{B}^T \mathbf{X}^T \mathbf{C} \mathbf{X} \mathbf{B}) = \mathbf{C}^T \mathbf{X} \mathbf{B} \mathbf{B}^T + \mathbf{C} \mathbf{X} \mathbf{B} \mathbf{B}^T.$$

**2.2. Examples and further discussions.** In this subsection, we give an example of polynomial regression on a synthetic data. We will sample various instances of the following random variable

$$Y = \sin(2\pi X) + \varepsilon, \quad (15)$$

where  $X \sim \text{Uniform}([0, 1])$  and  $\varepsilon \sim N(0, \sigma^2)$  with  $\sigma = 0.35$ , where  $X$  and  $\varepsilon$  are drawn independently.  $\varepsilon$  represents white noise, where the true signal should lie on the sine curve  $y = \sin(2\pi x)$ , which is unknown. We can control the signal-noise ratio through the parameter  $\sigma$ , which is the standard deviation of the Gaussian noise  $\varepsilon$ . For each  $N \in \{8, 100\}$ , we generate  $N$  i.i.d. samples of  $Y$  to form the train set, and another  $N$  i.i.d. samples of  $Y$  forms the test set. The goal is to fit the degree- $M$  polynomial regression model (3) to the train set and learn the optimal coefficients  $\hat{\mathbf{w}} \in \mathbb{R}^{M+1}$  by solving the least squares in (7), and use it to predict the test set via (10). Recall the closed-form solution for the optimal parameter  $\hat{\mathbf{w}}$  given in (9).

First consider  $N = 8$  samples with  $M = 1$  degree polynomial regression as shown in Figure 2 left. We are only given the  $N = 8$  training points represented by the blue dots. The red line represents the

<sup>2</sup> $\mathbf{X}^T \mathbf{X}$  is symmetric and positive semidefinite, and it is invertible iff the singular values of  $\mathbf{X}$  are all nonzero.

<sup>3</sup>The matrix  $\mathbf{X}^\dagger := (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$  is called the *Moore-Penrose pseudo-inverse* of  $\mathbf{X}$ . If  $\mathbf{X}$  is square and invertible, then  $\mathbf{X}^\dagger = \mathbf{X}^{-1} (\mathbf{X}^T)^{-1} \mathbf{X}^T = \mathbf{X}^{-1}$ . So the pseudo-inverse can be regarded as a generalization of matrix inverse for non-square matrices.

<sup>4</sup>Hint: Show that  $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$  is positive definite if  $\lambda > 0$ . Use the fact that the eigenvalues of a positive definite matrix  $\mathbf{A}$  has to be positive (why?), so  $\mathbf{Ay} \neq \mathbf{0}$  for any  $\mathbf{y}$  (why?) so  $\mathbf{A}$  is invertible (why?).

fitted polynomial, which is a linear function being a degree  $M = 1$  polynomial. Increasing the degree  $M$  of the polynomial model (and hence the degree of freedom of our model) to  $M = 4$ , we see the quartic polynomial in Figure 2 middle gives a much better fit to the sine curve. However, further increasing  $M$  to 9, the polynomial shown in Figure 2 right oscillates more wildly around the sine curve, giving a less favorable explanation of the true signal. While this may seem somewhat counterintuitive as the  $M = 9$  polynomial model entails the  $M = 4$  polynomial model, what happens here is an issue known as *overfitting*. Indeed, notice that the  $M = 9$  polynomial in Figure 2 right passes through every single train data point, while the  $M = 4$  polynomial in Figure 2 middle does not. Since the train set contains noise, our model should not be fitted too much to the observed data in the train set as it may learn ‘false patterns’ misled by the noise.

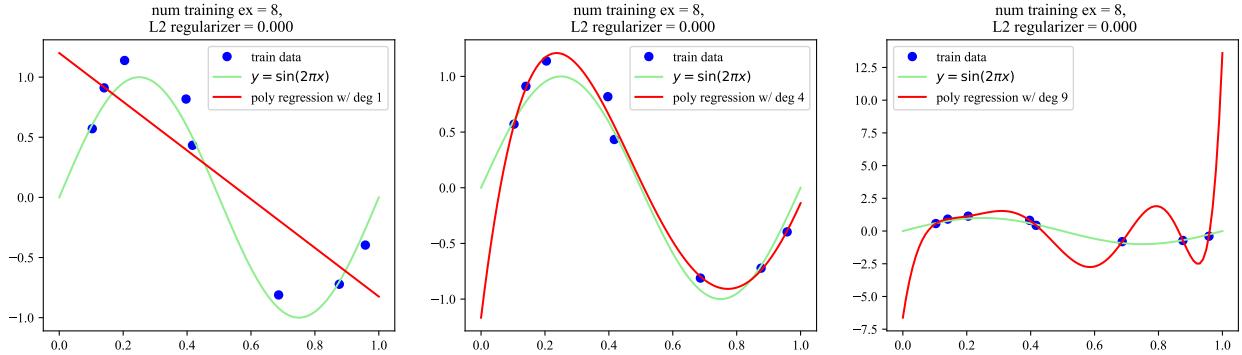


FIGURE 2. Example of polynomial regression fitting.  $N = 8$  samples for the train set and  $\lambda = 0$   $L_2$ -regularization parameter.

One way to overcome the issue of overfitting is to collect more training data. Indeed, if we use  $N = 100$  data points for the train set, the  $M = 9$  polynomial regression in Figure 3 right becomes more accurate to the true sine curve than the  $M = 4$  case in Figure 3 middle. Since the Gaussian noise in  $N = 100$  train samples are all independent, it will be extremely unlikely that many of the points lie on any given  $M = 9$  degree polynomial. Even though our synthetic data was built this way, this intuition is often useful in real-world data analysis tasks – noise in nature is close to Gaussian, and fitting to large train set learns ‘average patterns’, which should be close to the patterns in the original signal.

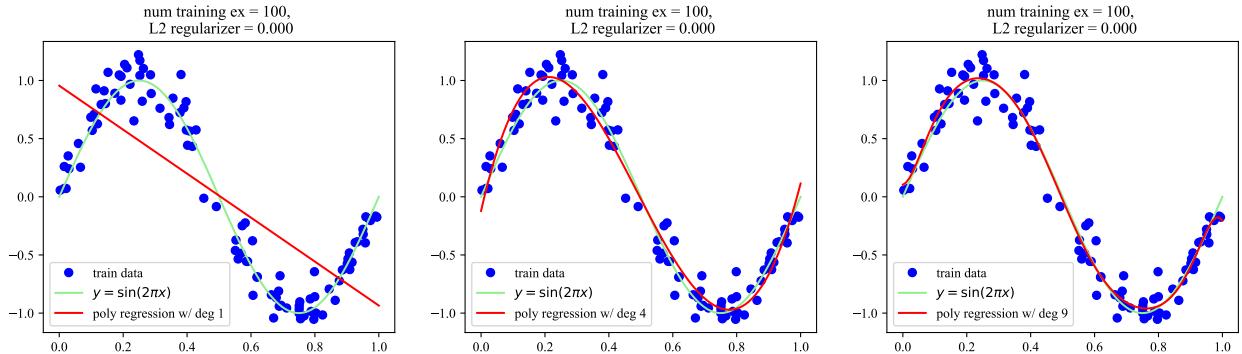


FIGURE 3. Example of polynomial regression fitting.  $N = 100$  samples for the train set and  $\lambda = 0$   $L_2$ -regularization parameter.

Another way to reduce overfitting is to use *regularization*. Namely, this is to put some penalization on the parameters of the model so that they cannot be too large. For example, consider the following

optimization problem

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^{M+1}} \|\mathbf{Y}_{\text{train}} - \mathbf{X}_{\text{train}} \mathbf{w}\|_F^2 + \lambda \|\mathbf{w}\|_F^2, \quad (16)$$

where we have added the ' $L_2$ -regularization term'  $\lambda \|\mathbf{w}\|_F^2$  to the squared loss  $\|\mathbf{Y}_{\text{train}} - \mathbf{X}_{\text{train}} \mathbf{w}\|_F^2$ . Here  $\lambda \geq 0$  is called the  $L_2$ -regularizer, which is a new hyperparameter that we introduced to the model (previously  $M$  was the only hyperparameter). Due to this additional penalization term, the coefficients of  $\hat{\mathbf{w}}$  in (16) cannot be too large. This technique is also called *shrinkage methods* in statistics as it shrinks the magnitude of the coefficients, and *weight decay* in the context of neural networks. The particular case of a quadratic regularizer as in (16) is called *ridge regression* [HK70]. If we instead use the  $L_1$ -regularization  $\lambda \|\mathbf{w}\|_1$ , then this is known as the (Lagrangian form of) the LASSO regression [Tib96], which better promotes 'sparsity' of  $\hat{\mathbf{w}}$ . An advantage of using the  $L_2$ -regularization is that we still have a closed-form solution for (16) by Exercise 1.2.1:

$$\hat{\mathbf{w}} = (\mathbf{X}_{\text{train}}^T \mathbf{X}_{\text{train}} + \lambda \mathbf{I})^{-1} \mathbf{X}_{\text{train}}^T \mathbf{Y}_{\text{train}}. \quad (17)$$

Indeed, using the  $L_2$ -regularizer  $\lambda = 0.001$  in the  $M = 9$  polynomial regression improves the performance for  $N = 8$  training data points, as shown in Figure 4 right.

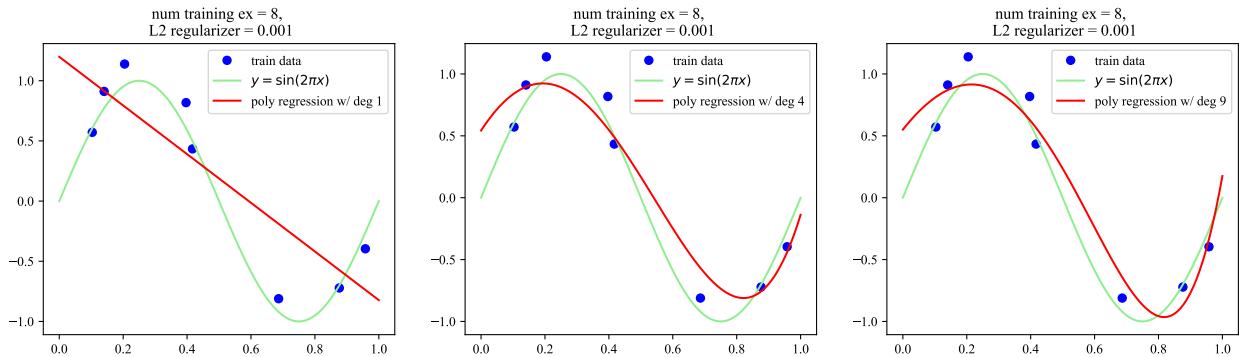


FIGURE 4. Example of polynomial regression fitting.  $N = 8$  samples for the train set and  $\lambda = 0.001$   $L_2$ -regularization parameter.

Lastly, we evaluate the performance of our learned polynomial model on the  $N$  test data points. We use the *mean squared error* (MSE)

$$\text{MSE}(\hat{\mathbf{w}}) := \frac{1}{N} \|\mathbf{Y}_{\text{test}} - \hat{\mathbf{Y}}_{\text{test}}\|_F^2, \quad (18)$$

where  $\hat{\mathbf{Y}}_{\text{test}} = \mathbf{X}_{\text{test}} \hat{\mathbf{w}}$  and  $\hat{\mathbf{w}}$  is given by (17). We use a similar MSE for the train set.

In Figure 5, we show the MSE for train and test sets against the degree  $M$  for three  $L_2$ -regularization parameters  $\lambda \in \{0, 0.001, 0.01\}$  for  $N \in \{8, 100\}$  data points both for the test and train sets. While the train error (in blue) mostly decreases as we increase the degree  $M$  in all cases, the test error starts to increase rapidly at  $M = 6$  with no regularization  $\lambda = 0$  (Figure 5 left). This indicates severe overfitting occurs past  $M = 6$ . The gap between the test and train error is also known as the *generalization error*<sup>5</sup>. In Figures 5 middle and right, we see the use of moderate  $L_2$ -regularizer  $\lambda = 0.001$  could significantly reduce overfitting (and hence the generalization error), but larger values of regularizer  $\lambda = 0.01$  hinders both the test and train errors. In general, the right value of the regularization parameter can be estimated by *cross-validation*. However, the performance of our model on hyperparameter becomes less sensitive when we have a large training set ( $N = 100$ ), as we see from Figure 6.

<sup>5</sup>We learn patterns from the observed test set. If the learned patterns capture the generic features of the underlying signal, then they should apply to the unobserved test set and yield a small test error; otherwise, then the pattern we learn does not 'generalize' beyond the test set.

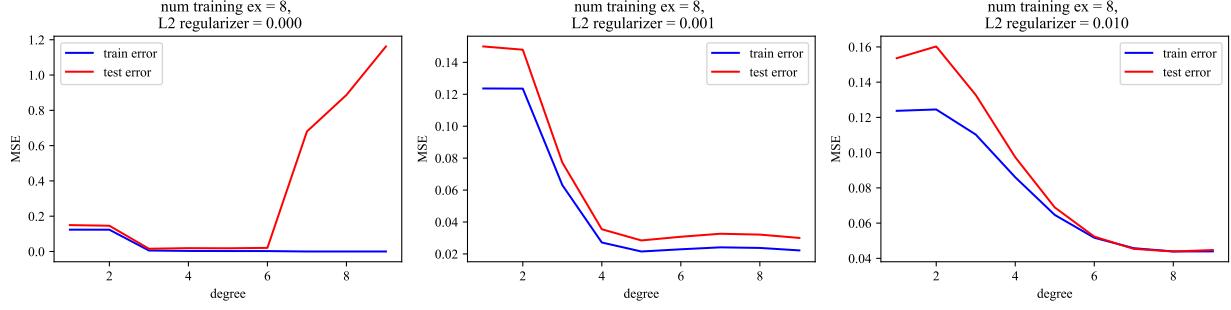


FIGURE 5. Plot of test and train error with respect to the degree  $M$  for three  $L_2$ -regularization parameters.  $N = 8$  points in both train and test sets.

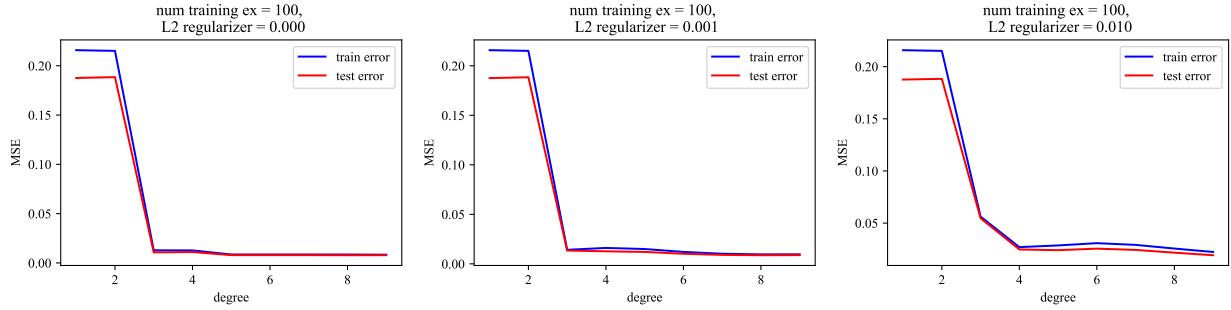


FIGURE 6. Plot of test and train error with respect to the degree  $M$  for three  $L_2$ -regularization parameters.  $N = 100$  points in both train and test sets.

**Exercise 1.2.3.** Using the [Jupyter notebook](#) provided in the course repository, reproduce the Figures 2 - 6, where the data are independently generated from the following random variable

$$Y = \cos(2\pi X) + \varepsilon, \quad (19)$$

where  $X \sim \text{Uniform}([0, 1])$  and  $\varepsilon \sim N(0, 0.16)$  are independent.

**2.3. Polynomial Regression by Maximum Likelihood.** In this subsection, we consider the polynomial regression with *Gaussian noise*. Namely, given an input  $x$ , the model predicts the output as a *random variable* (RV)  $\hat{y}(x; \mathbf{w}, \sigma) = \hat{y}(x; \mathbf{w}) + \varepsilon$ , where  $\hat{y}(x; \mathbf{w})$  is the polynomial function given in (3) and  $\varepsilon \sim N(0, \sigma^2)$  is an independent Gaussian noise with mean 0 and variance  $\sigma^2$ . Note that the predicted output  $\hat{y}(x; \mathbf{w}, \sigma)$  itself is a Gaussian RV with mean  $\hat{y}(x; \mathbf{w})$  and variance  $\sigma^2$  (see (10)):

$$\hat{y}(x; \mathbf{w}, \sigma) \sim N(\hat{y}(x; \mathbf{w}), \sigma^2), \quad \hat{y}(x; \mathbf{w}) = \sum_{k=0}^M w_k x^k = \boldsymbol{\phi}(x)^T \mathbf{w}. \quad (20)$$

This is a probabilistic model of polynomial regression with parameters  $\mathbf{w} = [w_0, \dots, w_M]^T$  for the coefficients of the unknown polynomial and  $\sigma$  for the standard deviation of the Gaussian noise.

Suppose we are given with a sequence of training examples  $[(x_1, y_1), \dots, (x_N, y_N)]$ . In order to fit the probabilistic polynomial regression (20) to the training dataset, we assume that the observed outputs  $y_1, \dots, y_N$  are realizations of the independent RVs  $\hat{y}_1(x_1; \mathbf{w}, \sigma), \dots, \hat{y}_N(x_N; \mathbf{w}, \sigma)$  satisfying (20). According to the method of *maximum likelihood* (see Section 6), we will have to find values of the parameters  $\mathbf{w}$

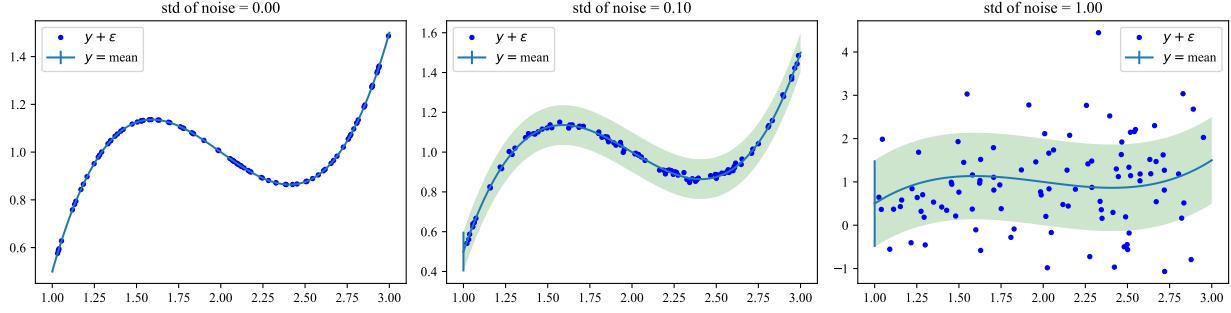


FIGURE 7. Plot of sample points of  $Y = -6 + (23/2)x - 6x^2 + x^3 + \varepsilon$ ,  $\varepsilon \sim N(0, \sigma^2)$ . Shaded regions represents  $\pm 1$  standard deviation of the noise.

and  $\sigma$  for which the joint likelihood  $L(y_1, \dots, y_N; \mathbf{w}, \sigma)$  of observing the tuple  $(y_1, \dots, y_N)$  from the tuple of independent RVs  $\hat{\mathbf{y}}_1(x_1; \mathbf{w}, \sigma), \dots, \hat{\mathbf{y}}_N(x_N; \mathbf{w}, \sigma)$  is maximized. Namely, the joint likelihood is given by

$$L(y_1, \dots, y_N; \mathbf{w}, \sigma) = \prod_{i=1}^n f_{\hat{\mathbf{y}}(x_i; \mathbf{w}, \sigma)}(y_i) \quad (21)$$

$$= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[ -\frac{(y_i - \hat{\mathbf{y}}(x_i; \mathbf{w}))^2}{2\sigma^2} \right] \quad (22)$$

$$= (2\pi\sigma^2)^{-N/2} \exp \left[ -\frac{1}{2\sigma^2} \sum_{k=1}^N (y_k - \hat{\mathbf{y}}(x_k; \mathbf{w}))^2 \right]. \quad (23)$$

A key observation to make is the following: Denoting  $\mathbf{Y}_{\text{train}} := [y_1, \dots, y_N]^T$ ,

$$\sum_{k=1}^N (y_k - \hat{\mathbf{y}}(x_k; \mathbf{w}))^2 = \|\mathbf{Y}_{\text{train}} - \mathbf{X}_{\text{train}}\mathbf{w}\|_F^2, \quad (24)$$

where  $\mathbf{X}_{\text{train}}$  is defined in (5) (see Exercise 1.2.4). Thus we can write the log likelihood function as

$$\log L(y_1, \dots, y_N; \mathbf{w}, \sigma) = -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|\mathbf{Y}_{\text{train}} - \mathbf{X}_{\text{train}}\mathbf{w}\|_F^2. \quad (25)$$

Noting that  $\log$  is strictly increasing function, we can find the *maximum likelihood estimate* (MLE) of the parameters  $\mathbf{w}, \sigma$  by the following optimization problem

$$(\hat{\mathbf{w}}, \hat{\sigma}) := \underset{\mathbf{w} \in \mathbb{R}^{M+1}, \sigma \geq 0}{\operatorname{argmax}} L(y_1, \dots, y_N; \mathbf{w}, \sigma) \quad (26)$$

$$= \underset{\mathbf{w} \in \mathbb{R}^{M+1}, \sigma \geq 0}{\operatorname{argmax}} \log L(y_1, \dots, y_N; \mathbf{w}, \sigma) \quad (27)$$

$$= \underset{\mathbf{w} \in \mathbb{R}^{M+1}, \sigma \geq 0}{\operatorname{argmax}} \left( -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|\mathbf{Y}_{\text{train}} - \mathbf{X}_{\text{train}}\mathbf{w}\|_F^2 \right). \quad (28)$$

Since the first term in the log likelihood function in (25) does not depend on  $\mathbf{w}$ , we see that  $\hat{\mathbf{w}}$  is the maximizer of  $-\frac{1}{2\sigma^2} \|\mathbf{Y}_{\text{train}} - \mathbf{X}_{\text{train}}\mathbf{w}\|_F^2$ . It follows that

$$\hat{\mathbf{w}} = \underset{\mathbf{w} \in \mathbb{R}^{M+1}}{\operatorname{argmin}} \|\mathbf{Y}_{\text{train}} - \mathbf{X}_{\text{train}}\mathbf{w}\|_F^2 = (\mathbf{X}_{\text{train}}^T \mathbf{X}_{\text{train}})^{-1} \mathbf{X}_{\text{train}}^T \mathbf{Y}_{\text{train}}, \quad (29)$$

where the second equality follows from Exercise 1.2.1. Note that this agrees with the optimal parameter  $\hat{\mathbf{w}}$  in (10) for the polynomial regression without noise.

On the other hand, in order to find  $\hat{\sigma}$ , we take the first and the second partial derivatives of (25) in  $\sigma$ :

$$\frac{\partial}{\partial \sigma} \log L(y_1, \dots, y_N; \mathbf{w}, \sigma) = -\frac{N}{\sigma} + \frac{1}{\sigma^3} \|\mathbf{Y}_{\text{train}} - \mathbf{X}_{\text{train}} \mathbf{w}\|_F^2, \quad (30)$$

$$\frac{\partial^2}{\partial \sigma^2} \log L(y_1, \dots, y_N; \mathbf{w}, \sigma) = -\frac{N}{\sigma^2} - \frac{3}{\sigma^4} \|\mathbf{Y}_{\text{train}} - \mathbf{X}_{\text{train}} \mathbf{w}\|_F^2 \leq 0. \quad (31)$$

According to the second derivative, the log likelihood is a concave function in  $\sigma$ . According to the first derivative, the log likelihood has a unique critical point. Since  $(\hat{w}, \hat{\sigma})$  jointly maximizes the log likelihood, we deduce

$$\hat{\sigma} = \frac{1}{N} \|\mathbf{Y}_{\text{train}} - \mathbf{X}_{\text{train}} \hat{\mathbf{w}}\|_F^2, \quad (32)$$

where  $\hat{\mathbf{w}}$  is given in (29).

Hence for any new input  $x$ , the optimal polynomial regression model will tell us to predict the corresponding output  $y$  as the following RV

$$\hat{Y}(x; \hat{\mathbf{w}}, \hat{\sigma}) \sim N(\phi(x)^T \hat{\mathbf{w}}, \hat{\sigma}^2), \quad (33)$$

where the optimal parameters  $\hat{\mathbf{w}}$  and  $\hat{\sigma}$  are given in (29) and (32), respectively.

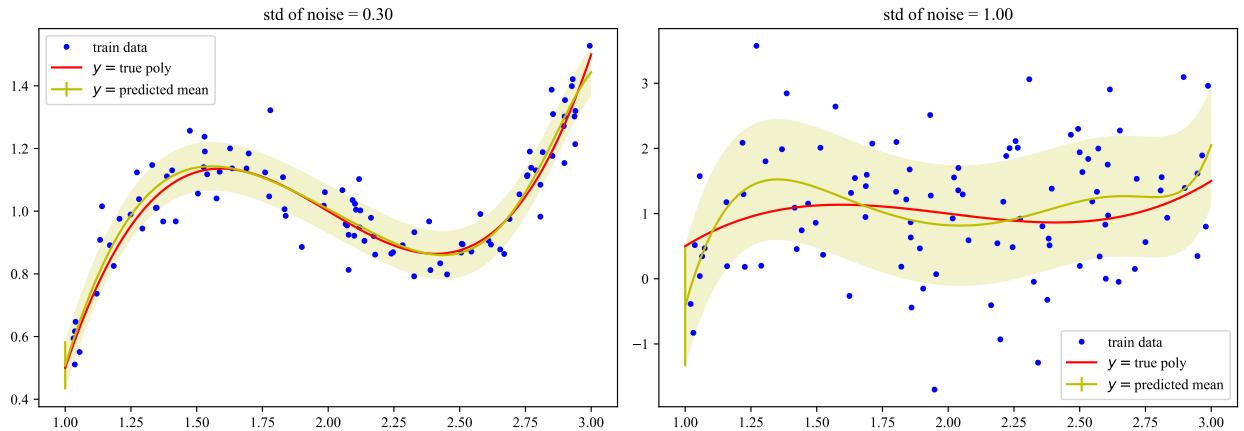


FIGURE 8. Plot of sample points of  $Y = -6 + (23/2)x - 6x^2 + x^3 + \varepsilon$ ,  $\varepsilon \sim N(0, \sigma^2)$  in blue for  $\sigma = 0.3$  (left) and  $\sigma = 1$  (right). Degree  $M = 9$  polynomial regression with Gaussian noise is fitted by maximum likelihood. The predicted mean is shown in the yellow curve with  $\pm 1$  predicted standard deviation ( $\hat{\sigma}$ ) shown in yellow shade.

**Exercise 1.2.4.** Fix  $\mathbf{w} = [w_0, w_1, \dots, w_M] \in \mathbb{R}^{M+1}$  and  $\sigma \geq 0$ . Let  $\hat{Y}_1, \dots, \hat{Y}_N$  be independent Gaussian RVs where  $\hat{Y}(x_i; \mathbf{w}, \sigma) \sim N(\phi(x)^T \mathbf{w}, \sigma^2)$  for  $i \in \{1, \dots, N\}$ , where  $\phi(x) = [1, x, \dots, x^M]^T$ . Show that the joint likelihood function for observing the values  $y_1, \dots, y_N$  is given by

$$L(y_1, \dots, y_N; \mathbf{w}, \sigma) = (2\pi\sigma^2)^{-N/2} \exp \left[ -\frac{1}{2\sigma^2} \|\mathbf{Y} - \mathbf{X}\mathbf{w}\|_F^2 \right], \quad (34)$$

where

$$\mathbf{Y} = [y_1, \dots, y_N]^T \in \mathbb{R}^N, \quad \mathbf{X} = \begin{bmatrix} 1 & x_1 & \cdots & x_1^M \\ \vdots & & & \\ 1 & x_2 & \cdots & x_N^M \end{bmatrix} \in \mathbb{R}^{N \times M}. \quad (35)$$

**2.4. Bayesian Polynomial Regression.** In this subsection, we take a Bayesian approach of the polynomial regression by regarding the  $\mathbf{w} = [w_0, \dots, w_M]^T \in \mathbb{R}^{M+1}$  as a (vector-valued) RV. For simplicity, we assume  $\sigma$  in (20) is known.

We assume that the components of  $\mathbf{w}$  are i.i.d. Gaussian RV with distribution  $N(0, \tau^2)$  for some  $\tau \geq 0$ . Equivalently,  $\mathbf{w}$  follows the multivariate Gaussian distribution  $N(\mathbf{0}, \tau\mathbf{I})$  with mean vector  $\mathbf{0} = [0, \dots, 0]^T \in \mathbb{R}^{M+1}$  and covariance matrix  $\tau\mathbf{I}$ , where  $\mathbf{I}$  is the  $(M+1) \times (M+1)$  identity matrix. Its probability distribution function (PDF) is given by

$$p(\mathbf{x}; \mathbf{0}, \tau\mathbf{I}) = (2\pi\tau^2)^{-(M+1)/2} \exp\left[-\frac{1}{2\tau^2}\mathbf{x}^T\mathbf{x}\right]. \quad (36)$$

See Example 3.2.2 for the general multivariate Gaussian distribution.

As briefly introduced in Section 7, Bayesian approach assumes parameters are random and continuously updates a hypothetical probability distribution on the parameters (a.k.a. *prior*) to the new distribution (a.k.a. *posterior*) with respect to the newly observed data according to Bayes' Theorem A.7.1. Namely, as in Subsection 2.3, suppose we are given with training examples  $[(x_1, y_1), \dots, (x_N, y_N)]$ . According to Bayes' Theorem, Exercise 1.2.4, and (36), the posterior distribution of the coefficients  $\mathbf{w}$  is given by the following proportionality relation

$$\overbrace{p(\mathbf{w}|\text{Data})}^{\text{posterior}} \propto \overbrace{L(\text{Data}|\mathbf{w})}^{\text{likelihood}} \overbrace{p(\mathbf{w})}^{\text{prior}} \quad (37)$$

$$= (2\pi\sigma^2)^{-N/2} \exp\left[-\frac{1}{2\sigma^2}\|\mathbf{Y} - \mathbf{X}\mathbf{w}\|_F^2\right] (2\pi\tau^2)^{-(M+1)/2} \exp\left[-\frac{1}{2\tau^2}\mathbf{w}^T\mathbf{w}\right] \quad (38)$$

$$= (2\pi\sigma^2)^{-N/2} (2\pi\tau^2)^{-(M+1)/2} \exp\left[-\frac{1}{2\sigma^2}\|\mathbf{Y} - \mathbf{X}\mathbf{w}\|_F^2 - \frac{1}{2\tau^2}\|\mathbf{w}\|_F^2\right]. \quad (39)$$

Later, we will show that the posterior distribution is in fact a multivariate Gaussian

$$p(\mathbf{w}|\text{Data}) = N(\mathbf{m}, \mathbf{S}), \quad (40)$$

where

$$\mathbf{m} := \sigma^{-2} \mathbf{S} \mathbf{X}_{\text{train}}^T \mathbf{Y}_{\text{train}} \in \mathbb{R}^{M+1}, \quad (41)$$

$$\mathbf{S}^{-1} := \tau^{-2} \mathbf{I} + \sigma^{-2} \mathbf{X}_{\text{train}}^T \mathbf{X}_{\text{train}} \in \mathbb{R}^{(M+1) \times (M+1)}. \quad (42)$$

Also, notice that the above posterior distribution is maximized when the exponent is maximized. Hence

$$\hat{\mathbf{w}}_{\text{MAP}} := \arg \max_{\mathbf{w} \in \mathbb{R}^{M+1}} p(\mathbf{w}; \text{Data}) = \arg \min_{\mathbf{w} \in \mathbb{R}^{M+1}} \left( \|\mathbf{Y} - \mathbf{X}\mathbf{w}\|_F^2 + \frac{\sigma^2}{\tau^2} \|\mathbf{w}\|_F^2 \right). \quad (43)$$

The quantity in the left hand side is called *maximum a posteriori* (MAP) estimate of  $\mathbf{w}$ . Note that the minimization problem in the right hand side is equivalent to (16) with the  $L_2$ -regularization parameter  $\lambda = (\sigma/\tau)^2$ .

Finally, for any new input  $x$ , the Bayesian polynomial regression model above will tell us to predict the corresponding output  $y$  as the following RV

$$\hat{Y}(x; \text{Data}, \sigma) \sim N(\boldsymbol{\phi}(x)^T \mathbf{w}_{\text{pos}}, \sigma^2), \quad (44)$$

where  $\mathbf{w}_{\text{pos}}$  is a *random* parameter following the posterior distribution given in (39). In fact, it can be shown that the distribution of  $\hat{Y}(x; \text{Data}, \sigma)$  is a (univariate) Gaussian:

$$\hat{Y}(x; \text{Data}, \sigma) \sim N(\boldsymbol{\phi}(x)^T \mathbf{m}, \sigma^2 + \boldsymbol{\phi}(x)^T \mathbf{S} \boldsymbol{\phi}(x)). \quad (45)$$

These claims will be justified later.

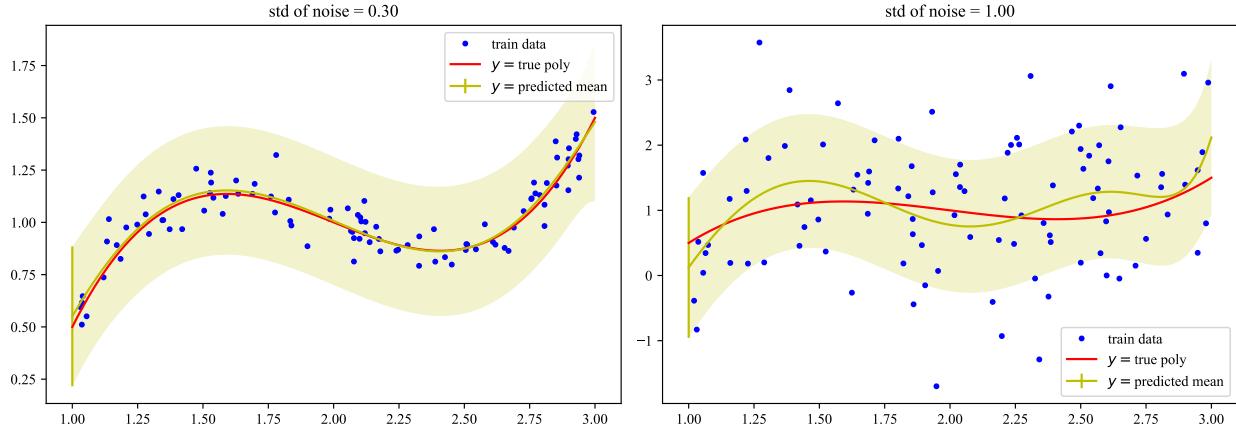


FIGURE 9. Plot of sample points of  $Y = -6 + (23/2)x - 6x^2 + x^3 + \varepsilon$ ,  $\varepsilon \sim N(0, \sigma^2)$  in blue for  $\sigma = 0.3$  (left) and  $\sigma = 1$  (right). Degree  $M = 9$  Bayesian polynomial regression with prior  $\mathbf{w} \sim N(\mathbf{0}, \tau^2 \mathbf{I})$ ,  $\tau = 10$  is fitted. The predicted mean is shown in the yellow curve with  $\pm 1$  predicted standard deviation ( $\hat{\sigma}$ ) shown in yellow shade.

**Exercise 1.2.5.** Using the [Jupyter notebook](#) provided in the course repository, reproduce the Figures 8 - 9, where the data are independently generated from the following random variable

$$Y = \cos(2\pi X) + \varepsilon, \quad (46)$$

where  $X \sim \text{Uniform}([0, 1])$  and  $\varepsilon \sim N(0, 0.16)$  are independent. Compare the results of the maximum likelihood and the Bayesian polynomial regression.

### 3. Decision Theory

Suppose we have an input vector  $\mathbf{x} \in \mathcal{X}$  together with a corresponding output vector  $\mathbf{y} \in \mathcal{Y}$ . Our goal is to predict the output  $\mathbf{y}$  given a new value for  $\mathbf{x}$ , which will take a continuous value for regression problems and discrete value (class label) for classification problems. More formally, a *decision algorithm* is a function  $\mathcal{A} : \mathcal{X} \rightarrow \mathcal{Y}$  that maps a new input vector  $\mathbf{x} \in \mathcal{X}$  to a predicted output  $\mathbf{y} \in \mathcal{Y}$ .

Method of learning decision algorithms can be roughly divided into ‘discriminative’ and ‘generative’ ones. *Discriminative* methods learn decision algorithms directly from training dataset without any probabilistic modeling of the relationship between the input  $\mathbf{x}$  and the output  $\mathbf{y}$ . For instance, the polynomial regression model  $\mathbf{x} \mapsto \hat{y}(\mathbf{x}; \mathbf{w}) = \phi(x)^T \mathbf{w}$  in (10) is a classical example of discriminative decision algorithms for the case of regression problems. On the other hand, a classical example of discriminative binary classification algorithm is the following *perceptron algorithm*:

$$\mathbf{x} \longmapsto f(\phi(\mathbf{x})^T \mathbf{w}) \in \{+1, -1\}, \quad (47)$$

where  $\phi(\mathbf{x})$  is a *feature vector* of the input  $\mathbf{x}$  (e.g.,  $\phi(x) = [1, x, \dots, x^M]^T$  in the case of polynomial regression) and  $f$  is the *activation function* of the form

$$f(x) = \mathbf{1}(x \geq \theta) - \mathbf{1}(x < \theta) = \begin{cases} +1 & x \geq \theta \\ 1 & x < \theta \end{cases}, \quad (48)$$

where  $\theta \geq 0$  is a threshold parameter.

In the rest of this section, we take a closer look at the generative methods of learning decision algorithms. This is roughly a two-step process involving the ‘inference problem’ and the ‘decision problem’.

**1. Inference problem:** We model the unknown input and output as the random vector  $(\mathbf{X}, \mathbf{Y})$  following the joint distribution  $(\mathbf{x}, \mathbf{y}) \mapsto p(\mathbf{x}, \mathbf{y})$ :

$$\mathbb{P}(\mathbf{X} = \mathbf{x}, \mathbf{Y} = \mathbf{y}) = p(\mathbf{x}, \mathbf{y}). \quad (49)$$

This describes our complete understanding on the pattern between  $\mathbf{x}$  and  $\mathbf{y}$ , where randomness represents the uncertainty associated to these variables. Conditioning on the input value  $\mathbf{x}$ , we can also deduce that the output  $\mathbf{Y}$  given the observed input  $\mathbf{x}$  has the following conditional distribution

$$p(\mathbf{y}|\mathbf{x}) := \mathbb{P}(\mathbf{Y} = \mathbf{y} | \mathbf{X} = \mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x})}, \quad (50)$$

where the marginal distribution of  $\mathbf{X}$  is given by

$$p(\mathbf{x}) = \int_{\mathcal{Y}} p(\mathbf{x}, \mathbf{y}) d\mathbf{y}. \quad (51)$$

It is possible that one only compute the conditional distribution  $p(\mathbf{y}|\mathbf{x})$  in the inference step and proceed to the decision problem (e.g., logistic regression).

**2. Decision problem:** Given a conditional probability distribution  $\mathbf{y} \mapsto p(\mathbf{y}|\mathbf{x})$  and given a new value of input  $\mathbf{x}$ , decide (predict) the value of corresponding output  $\mathbf{y}$ . A general principle to find an optimal decision algorithm using the conditional distribution is the following (*conditional expected loss minimization*) (ELM):

$$\mathbf{x} \mapsto \mathcal{A}(\mathbf{x}) := \arg \min_{\hat{\mathbf{y}} \in \mathcal{Y}} \mathbb{E}[\ell(\mathbf{Y}, \hat{\mathbf{y}}) | \mathbf{X} = \mathbf{x}], \quad (52)$$

where  $(\mathbf{y}, \hat{\mathbf{y}}) \mapsto \ell(\mathbf{y}, \hat{\mathbf{y}}) \in \mathbb{R}$  is a chosen *loss function* that measures the ‘goodness’ of our prediction  $\hat{\mathbf{y}}$  with respect to the true output  $\mathbf{y}$ .

Below are some instances of decision algorithms found by ELM:

*Classification:* The prediction  $\hat{\mathbf{y}}$  is the maximizer of the conditional distribution  $p(\mathbf{y}|\mathbf{x})$ :

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y} \in \mathcal{Y}} p(\mathbf{y}|\mathbf{x}) = \text{‘most probable output given } \mathbf{x}\text{’}. \quad (53)$$

*Regression:* The prediction  $\hat{\mathbf{y}}$  is the conditional expectation of the random output  $\mathbf{Y}$  given  $\mathbf{x}$ :

$$\hat{\mathbf{y}} = \mathbb{E}[\mathbf{Y} | \mathbf{X} = \mathbf{x}] = \text{‘expected output given } \mathbf{x}\text{’}. \quad (54)$$

The inference problem is more demanding both theoretically and computationally than solving the decision problem. In practical applications, we must often make a specific prediction for the value of  $\mathbf{y}$  in order to take a suitable action based on our current best knowledge of the problem.

**Example 1.3.1** (Detecting cancer from x-ray images). Consider the situation where the input is an *x-ray* image of a new patient, and we would like to determine the output  $\mathbf{y}$  between two classes *normal* and *cancer*. Our decision algorithm  $\mathcal{A}$  maps the patient’s x-ray image  $\mathbf{x}$  to a class label  $\mathbf{y} \in \{\text{normal}, \text{cancer}\}$ . There are four possible outcomes depending on whether the patient has cancer or not and how she has been diagnosed by our algorithm. Since misclassifying a patient who has cancer may be far more detrimental than misclassifying a patient who does not have cancer, it would be reasonable to consider a loss function  $\ell$  of the following form:

$$\ell(\text{true}, \text{prediction}) = \begin{array}{ccc} \text{true} \backslash \text{prediction} & \text{normal} & \text{cancer} \\ \text{normal} & 0 & 1 \\ \text{cancer} & 1000 & 0 \end{array}. \quad (55)$$

Then the decision algorithm  $\mathcal{A}$  that minimizes the corresponding expected loss would be more aggressive in looking for cues of cancer and classifying as *cancer* in order to avoid the huge loss incurred when missing patient with cancer. ▲

**Example 1.3.2** (Maxizing conditional probabilities for classification). Consider the case of classification problem, where the output variable  $\mathbf{Y}$  is a discrete random variable (taking values of class labels). Consider the following ‘indicator loss function’

$$\ell(\mathbf{y}, \hat{\mathbf{y}}) = \mathbf{1}(\mathbf{y} \neq \hat{\mathbf{y}}). \quad (56)$$

Namely, we gain a unit loss whenever our prediction  $\hat{\mathbf{y}}$  differs from the true output  $\mathbf{y}$ . Then note that

$$\mathbb{E}[\ell(\mathbf{Y}, \hat{\mathbf{y}}) | \mathbf{X} = \mathbf{x}] = \mathbb{E}[\mathbf{1}(\mathbf{y} \neq \hat{\mathbf{y}}) | \mathbf{X} = \mathbf{x}] = \mathbb{P}(\mathbf{Y} \neq \hat{\mathbf{y}} | \mathbf{X} = \mathbf{x}) = 1 - \mathbb{P}(\mathbf{Y} = \hat{\mathbf{y}} | \mathbf{X} = \mathbf{x}) = 1 - p(\hat{\mathbf{y}} | \mathbf{x}). \quad (57)$$

Thus, the decision algorithm  $\mathcal{A}$  given by the ELM (52) in this case reduces to

$$\mathbf{x} \mapsto \mathcal{A}(\mathbf{x}) = \arg \min_{\hat{\mathbf{y}} \in \mathcal{Y}} \mathbb{E}[\ell(\mathbf{Y}, \hat{\mathbf{y}}) | \mathbf{X} = \mathbf{x}] = \arg \max_{\hat{\mathbf{y}} \in \mathcal{Y}} p(\hat{\mathbf{y}} | \mathbf{x}). \quad (58)$$

In words, we predict the true class label  $\mathbf{Y}$  given  $\mathbf{x}$  as the most probable class label given  $\mathbf{x}$ . ▲

**Exercise 1.3.3** (Minimizing squared loss for regression). Consider the case of regression problem, where the output variable  $\mathbf{Y} \in \mathbb{R}^d$  is a continuous  $d$ -dimensional random vector (e.g., stock prices in a portfolio tomorrow). Consider the squared loss:

$$\ell(\mathbf{y}, \hat{\mathbf{y}}) = \|\mathbf{y} - \hat{\mathbf{y}}\|_F^2 = (\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}}). \quad (59)$$

In this exercise, we will show that the decision algorithm  $\mathcal{A}$  given by the ELM (52) in this case reduces to

$$\mathbf{x} \mapsto \mathcal{A}(\mathbf{x}) = \arg \min_{\hat{\mathbf{y}} \in \mathcal{Y}} \mathbb{E}[\ell(\mathbf{Y}, \hat{\mathbf{y}}) | \mathbf{X} = \mathbf{x}] = \mathbb{E}[\mathbf{Y} | \mathbf{X} = \mathbf{x}]. \quad (60)$$

(i) Write  $\mathbf{Y} = [Y_1, Y_2, \dots, Y_d]^T$ . Show that

$$\mathbb{E}\left[\|\mathbf{Y}\|_F^2 \mid \mathbf{X} = \mathbf{x}\right] - \|\mathbb{E}[\mathbf{Y} | \mathbf{X} = \mathbf{x}]\|_F^2 = \sum_{i=1}^d \text{Var}(Y_i | \mathbf{X} = \mathbf{x}). \quad (61)$$

(ii) Justify each equality in the following ‘completing squares’ computation:

$$\mathbb{E}\left[\|\mathbf{Y} - \mathcal{A}(\mathbf{x})\|_F^2 \mid \mathbf{X} = \mathbf{x}\right] = \mathbb{E}\left[\mathbf{Y}^T \mathbf{Y} - 2\mathcal{A}(\mathbf{x})^T \mathbf{Y} + \mathcal{A}(\mathbf{x})^T \mathcal{A}(\mathbf{x}) \mid \mathbf{X} = \mathbf{x}\right] \quad (62)$$

$$= \mathbb{E}\left[\|\mathbf{Y}\|_F^2 \mid \mathbf{X} = \mathbf{x}\right] - 2\mathcal{A}(\mathbf{x})^T \mathbb{E}[\mathbf{Y} | \mathbf{X} = \mathbf{x}] + \|\mathcal{A}(\mathbf{x})\|_F^2 \quad (63)$$

$$= \|\mathcal{A}(\mathbf{x}) - \mathbb{E}[\mathbf{Y} | \mathbf{X} = \mathbf{x}]\|_F^2 + \sum_{i=1}^d \text{Var}(Y_i | \mathbf{X} = \mathbf{x}). \quad (64)$$

(iii) From (ii), deduce that the expected squared loss  $\mathbb{E}[\ell(\mathbf{Y}, \hat{\mathbf{y}}) | \mathbf{X} = \mathbf{x}]$  of the decision algorithm  $\mathbf{x} \mapsto \mathcal{A}(\mathbf{x})$  is minimized when  $\mathcal{A}(\mathbf{x}) = \mathbb{E}[\mathbf{Y} | \mathbf{X} = \mathbf{x}]$  and the global minimum equals  $\sum_{i=1}^d \text{Var}(Y_i | \mathbf{X} = \mathbf{x})$ .

In the following proposition, we show that minimizing ‘total expected loss’ and ‘conditional expected loss’ give the same decision algorithm.

**Proposition 1.3.4** (Total expected loss minimization). *Consider the following decision algorithm  $\mathcal{B} : \mathcal{X} \rightarrow \mathcal{Y}$  given by the following total expected loss minimization:*

$$\mathcal{B} = \arg \min_{f: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}[\ell(\mathbf{Y}, f(\mathbf{X}))], \quad (65)$$

where  $(\mathbf{y}, \hat{\mathbf{y}}) \mapsto \ell(\mathbf{y}, \hat{\mathbf{y}})$  is a chosen loss function. Let  $\mathcal{A}$  denote the decision algorithm given by the (conditional) expected loss minimization:

$$\mathbf{x} \mapsto \mathcal{A}(\mathbf{x}) = \arg \min_{\hat{\mathbf{y}} \in \mathcal{Y}} \mathbb{E}[\ell(\mathbf{Y}, \hat{\mathbf{y}}) | \mathbf{X} = \mathbf{x}]. \quad (66)$$

Then it holds that  $\mathcal{B} = \mathcal{A}$ .

PROOF. It suffices to show that  $\mathcal{A}$  achieves the minimum total expected loss:

$$\mathbb{E}[\ell(\mathbf{Y}, \mathcal{A}(\mathbf{X}))] = \min_{f: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}[\ell(\mathbf{Y}, f(\mathbf{X}))] = \mathbb{E}[\ell(\mathbf{Y}, \mathcal{B}(\mathbf{X}))], \quad (67)$$

where the second equality follows from the definition of  $\mathcal{B}$ . First, since  $\mathcal{B}$  minimizes the total expected loss and  $\mathcal{A}$  is one of the possible decision algorithms,

$$\mathbb{E}[\ell(\mathbf{Y}, \mathcal{A}(\mathbf{X}))] \geq \mathbb{E}[\ell(\mathbf{Y}, \mathcal{B}(\mathbf{X}))]. \quad (68)$$

On the other hand, by the minimality of  $\mathcal{A}$ , we have for each  $\mathbf{x} \in \mathcal{X}$ ,

$$\mathbb{E}_{\mathbf{Y}}[\ell(\mathbf{Y}, \mathcal{B}(\mathbf{X})) | \mathbf{X} = \mathbf{x}] \geq \mathbb{E}[\ell(\mathbf{Y}, \mathcal{A}(\mathbf{x})) | \mathbf{X} = \mathbf{x}] \quad (69)$$

By taking expectation over  $\mathbf{X}$  and using iterated expectation,

$$\mathbb{E}[\ell(\mathbf{Y}, \mathcal{B}(\mathbf{X}))] = \mathbb{E}_{\mathbf{X}}[\mathbb{E}_{\mathbf{Y}}[\ell(\mathbf{Y}, \mathcal{B}(\mathbf{X})) | \mathbf{X}]] \geq \mathbb{E}_{\mathbf{X}}[\mathbb{E}_{\mathbf{Y}}[\ell(\mathbf{Y}, \mathcal{A}(\mathbf{X})) | \mathbf{X}]] = \mathbb{E}[\ell(\mathbf{Y}, \mathcal{A}(\mathbf{X}))]. \quad (70)$$

Hence by combining (68) and (70), we conclude (67), as desired.  $\square$

## CHAPTER 2

# Algorithms for Classification

## 1. Introduction

A particular type of pattern recognition problem that machine learning methods have made significant progress over the past decades is *classification problems*. Here are some motivating examples:

1. Classify patients' chest x-ray images into classes 'normal' and 'pneumonia'.
2. Classify dog and cat images into classes 'dog' and 'cat'.
3. Classify an image of handwritten digit into the correct digit. (see Figure 10)
4. Classify incoming emails into 'spam' and 'not spam'.
5. Classify given user behavior on a social network into 'churn' or 'human'

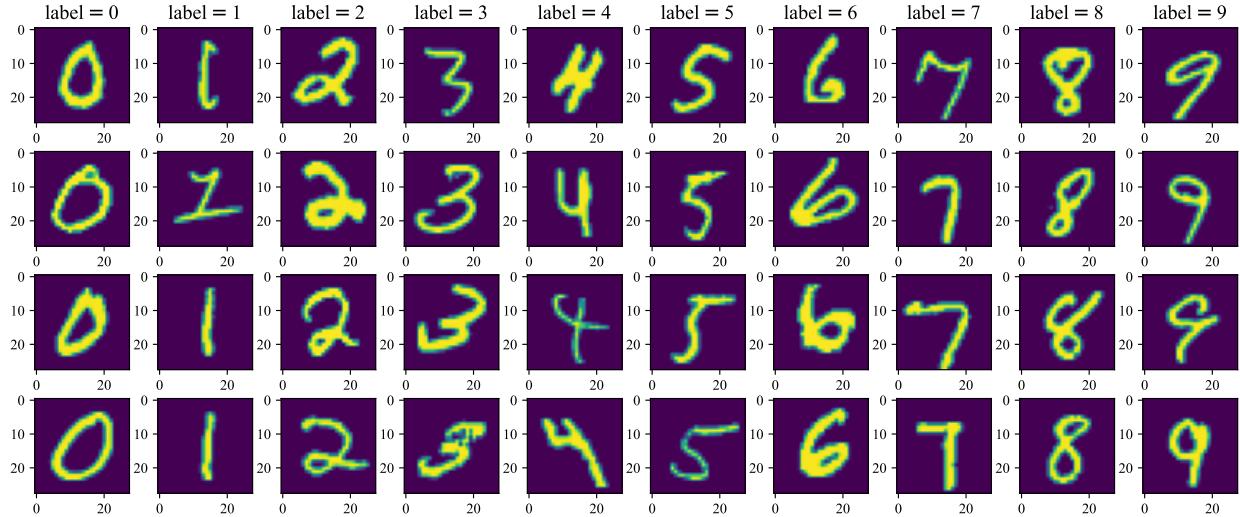


FIGURE 10. Examples of [MNIST](#) dataset of handwritten digits by Yann LeCun, Corinna Cortes, and Christopher J.C. Burges. There are total 70000 images of size  $28 \times 28$  pixels.

In general, for classification problem, we seek to learn a decision algorithm  $\mathcal{A} : \mathcal{X} \rightarrow \mathcal{Y} = \{\mathcal{C}_1, \dots, \mathcal{C}_K\}$ , which maps each object  $\mathbf{x} \in \mathcal{X}$  into one of  $K$  possible classes  $\mathcal{C}_1, \dots, \mathcal{C}_K$ . This problem is called a *binary classification* if  $K = 2$  and a *multi-class classification* if  $K \geq 3$ . For instance, the MNIST dataset consists of total 70000 black-and-white images of size  $28 \times 28$  pixels. Viewing each image as a matrix in  $\mathbb{R}^{28 \times 28}$ , the classification problem for the MNIST dataset  $\mathcal{X}_{\text{MNIST}}$  can be thought of as a multi-class classification problem where one wants to learn a decision algorithm  $\mathcal{A} : \mathbb{R}^{28 \times 28} \supseteq \mathcal{X}_{\text{MNIST}} \rightarrow \{0, 1, 2, \dots, 9\}$ . Equivalently, by vectorizing each image into a  $784 = 28^2$  dimensional vector, this is also to learn a decision algorithm  $\mathcal{A} : \subseteq \mathbb{R}^{784} \supseteq \mathcal{X}_{\text{MNIST}} \rightarrow \{0, 1, 2, \dots, 9\}$ .

Some of the popular classification algorithms are: Logistic Regression, k-Nearest Neighbors, Decision Trees, Support Vector Machine, and Naive Bayes. We will take a closer look at some of the algorithms in this chapter.

## 2. Logistic Regression

**2.1. The model formulation.** We begin our discussion for the simplest case of binary classification. Without loss of generality, we can assume that the two classes are 0 and 1. Say we have training examples  $(\phi(\mathbf{x}_1), y_1), \dots, (\phi(\mathbf{x}_N), y_N)$ , where

- $\mathbf{x}_1, \dots, \mathbf{x}_N$ : Input data (e.g., collection of all medical records of each patient)
- $\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_N) \in \mathbb{R}^p$ : Features (e.g., some useful (derived) information for each patient)
- $y_1, \dots, y_N \in \{0, 1\}$ : Binary class labels (e.g., ‘normal’ or ‘pneumonia’).

The basic idea of logistic regression is to model the output  $y$  as a Bernoulli (0-1) random variable  $Y \sim \text{Bernoulli}(p)$  with success probability  $p$  that depends on the observed feature  $\phi(\mathbf{x})$ . More specifically,

$$Y | \phi(\mathbf{x}) \sim \text{Bernoulli}(p), \quad p = \sigma(\phi(\mathbf{x})^T \mathbf{w}) := \frac{\exp(\phi(\mathbf{x})^T \mathbf{w})}{1 + \exp(\phi(\mathbf{x})^T \mathbf{w})} \quad (71)$$

Here,  $\sigma : x \mapsto e^x/(1 + e^x) = 1/(1 + e^{-x})$  is called a *sigmoid* (or *logistic*) function, and  $\mathbf{w} \in \mathbb{R}^p$  is a  $p$ -dimensional parameter vector. Notice that the sigmoid function maps the real line onto the interval  $[0, 1]$  (see Figure 11). Hence one can think of it as a function that ‘links’ real scalar with an associated probability. So in logistic regression, we model the ‘success probability’  $p$  given the feature  $\phi(\mathbf{x})$  as the dot product  $\phi(\mathbf{x})^T \mathbf{w} \in \mathbb{R}$  fed into the sigmoid function  $\sigma$  in order to make it a probability  $\in [0, 1]$ .

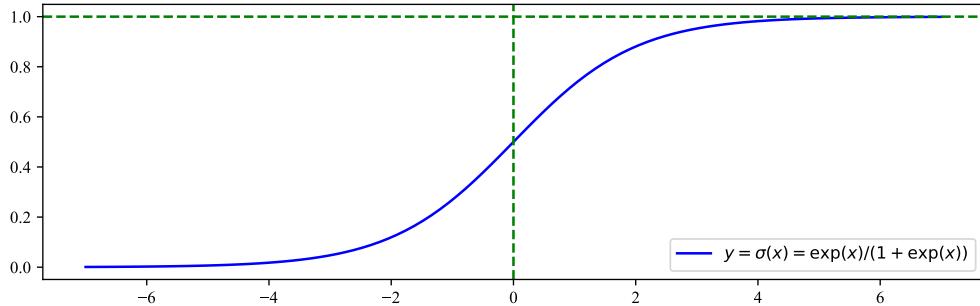


FIGURE 11. Plot of sigmoid function.

Next, let’s discuss how to train logistic regression model. Since it is a probabilistic model, we will use the method of maximum likelihood. For that, we compute the joint likelihood of observing the output values  $(y_1, \dots, y_N)$  from the tuple of independent Bernoulli variables  $(Y_1, \dots, Y_N)$  with success probabilities  $p_1 := \sigma(\phi(\mathbf{x}_1)^T \mathbf{w}), \dots, p_N := \sigma(\phi(\mathbf{x}_N)^T \mathbf{w})$ :

$$L(y_1, \dots, y_N; \mathbf{w}) = \mathbb{P}(Y_1 = y_1, \dots, Y_N = y_N; \mathbf{w}) = \prod_{i=1}^N p_i^{y_i} (1 - p_i)^{1-y_i}. \quad (72)$$

Hence the optimal parameter  $\hat{\mathbf{w}} \in \mathbb{R}^p$  is given by

$$\hat{\mathbf{w}} = \underset{\mathbf{w} \in \mathbb{R}^p}{\operatorname{argmax}} \log L(y_1, \dots, y_N; \mathbf{w}) \quad (73)$$

$$= \underset{\mathbf{w} \in \mathbb{R}^p}{\operatorname{argmax}} \sum_{i=1}^N y_i \log p_i + (1 - y_i) \log(1 - p_i) \quad (74)$$

$$= \underset{\mathbf{w} \in \mathbb{R}^p}{\operatorname{argmin}} \sum_{i=1}^N y_i \log(p_i^{-1} - 1) - \log(1 - p_i) \quad (75)$$

$$= \underset{\mathbf{w} \in \mathbb{R}^p}{\operatorname{argmin}} \sum_{i=1}^N -y_i \phi(\mathbf{x}_i)^T \mathbf{w} + \log(1 + \exp(\phi(\mathbf{x}_i)^T \mathbf{w})), \quad (76)$$

where we have used the following observations

$$p_i^{-1} - 1 = \exp(-\boldsymbol{\phi}(\mathbf{x}_i)^T \mathbf{w}), \quad (1 - p_i)^{-1} = 1 + \exp(\boldsymbol{\phi}(\mathbf{x}_i)^T \mathbf{w}). \quad (77)$$

Using a matrix formulation, we have

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^p} \left[ \ell_{LR}(\mathbf{w}) := \left( \sum_{i=1}^N \log(1 + \exp(\boldsymbol{\phi}(\mathbf{x}_i)^T \mathbf{w})) \right) - \mathbf{Y}^T \boldsymbol{\Phi}^T \mathbf{w} \right], \quad (78)$$

where

$$\mathbf{Y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} \in \{0, 1\}^N, \quad \boldsymbol{\Phi} = \begin{bmatrix} \boldsymbol{\phi}(\mathbf{x}_1) & \cdots & \boldsymbol{\phi}(\mathbf{x}_N) \end{bmatrix} \in \mathbb{R}^{p \times N}. \quad (79)$$

Here  $\boldsymbol{\Phi}$  is also called the *design matrix*. The loss function  $\ell_{LR}(\mathbf{w})$  in (78) is called the *logistic regression loss*.

**2.2. Minimizing the logistic regression loss.** In the previous section, we have formulated the logistic regression model and derived an optimization problem (78) for its training. Namely, we need to minimize the logistic regression loss  $\ell_{LR}$  defined in (78) in order to find the optimal (maximum likelihood) parameter  $\hat{\mathbf{w}}$ . Recall that in the case of polynomial regression, the quadratic optimization problem (7) had a closed-form solution for the optimal  $\hat{\mathbf{w}}$ . In the case of the logistic regression, we will observe that there is no closed-form solution for the minimizer, but it is still a *convex optimization problem*, which can be solved effectively by an iterative algorithm (e.g., gradient descent, Newton-Ralphson).

The following exercise provides the gradient, the Hessian, and convexity of the logistic regression loss  $\ell_{LR}$ .

**Exercise 2.2.1** (Convexity of the logistic regression loss). Let  $\ell_{LR}(\mathbf{w})$  denote the logistic regression loss defined in (78). Let  $\mathbf{Q} = \mathbf{Q}(\mathbf{w}) = [p_1, \dots, p_N]^T \in \mathbb{R}^N$  where  $p_i = \sigma(\boldsymbol{\phi}(\mathbf{x}_i)^T \mathbf{w})$  with  $\sigma(x) = e^x / (1 + e^x)$ .

- (i) Show that  $\frac{\partial \sigma(x)}{x} = \sigma(x)(1 - \sigma(x))$ .
- (ii) Show that

$$\frac{\partial \mathbf{Q}^T}{\partial \mathbf{w}} = \begin{bmatrix} \frac{\partial}{\partial w_1} \\ \vdots \\ \frac{\partial}{\partial w_p} \end{bmatrix} [p_1 \quad \cdots \quad p_N] = \begin{bmatrix} p_1(1 - p_1)\boldsymbol{\phi}(\mathbf{x}_1) & \cdots & p_N(1 - p_N)\boldsymbol{\phi}(\mathbf{x}_N) \end{bmatrix} = \boldsymbol{\Phi} \mathbf{D}, \quad (80)$$

where  $\mathbf{D}$  is the  $N \times N$  diagonal matrix with diagonal entries  $\mathbf{D}_{ii} = p_i(1 - p_i)$ .

- (iii) Show that

$$\nabla_{\mathbf{w}} \ell_{LR}(\mathbf{w}) = \boldsymbol{\Phi}(\mathbf{Q} - \mathbf{Y}) \in \mathbb{R}^p, \quad \mathbf{H} := \frac{\partial^2 \ell_{LR}(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^T} = \nabla_{\mathbf{w}} \mathbf{Q}^T \boldsymbol{\Phi}^T = \boldsymbol{\Phi} \mathbf{D} \boldsymbol{\Phi}^T \in \mathbb{R}^{p \times p}. \quad (81)$$

- (iv) Argue that the Hessian  $\mathbf{H}$  in (iii) is positive semi-definite. Conclude that the logistic regression loss  $\ell_{LR}$  is convex. (Hint:  $\mathbf{D}$  is a diagonal matrix with positive diagonal entries.)

Now, we will introduce two convex optimization algorithms for logistic regression. The algorithms are iterative in nature, meaning that they apply the same update rule until a certain convergence criterion is satisfied. The first algorithm belongs to the class of *gradient descent*:

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta_t \nabla_{\mathbf{w}} \ell_{LR}(\mathbf{w}_t) = \mathbf{w}_t - \eta_t \boldsymbol{\Phi}(\mathbf{Q}(\mathbf{w}_t) - \mathbf{Y}), \quad (82)$$

where  $(\eta_t)_{t \geq 0}$  is a sequence of positive numbers called *step sizes*. Recall that the gradient  $\nabla_{\mathbf{w}} \ell_{LR}$  is the direction of change in  $\mathbf{w}$  that increases the loss function  $\ell_{LR}$  most rapidly. Since we want to minimize the loss function, we would want to incrementally change the current iterate  $\mathbf{w}_t$  to the opposite direction of the gradient. The remaining degree of freedom is the step size that tells us how far we follow such

direction  $-\nabla_{\mathbf{w}} \ell_{LR}(\mathbf{w}_t)$  of the steepest descent<sup>1</sup>. A particular choice that is guaranteed to converge to a global minimum for convex optimization problems is  $\eta_t = \gamma(\log t)/\sqrt{t}$ , where  $\gamma > 0$  is a constant<sup>2</sup>.

The second algorithm also belongs to the class of gradient descent but uses the Hessian  $\mathbf{H} = \nabla^2 \ell_{LR}(\mathbf{w}_t)$ , which is called the *Newton-Ralphson method*. It takes the following form:

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \mathbf{H}^{-1} \nabla_{\mathbf{w}} \ell_{LR}(\mathbf{w}_t). \quad (83)$$

One can think of the above algorithm as a vresion of gradient descent where we used the Hessian inverse  $\mathbf{H}^{-1}$  as a ‘matrix-valued step size’<sup>3</sup>. The use of second-order derivatives often accelerates the speed of convergence but suffers from higher computational cost per iteration. In our current case of logistic regression, the Newton-Ralphson algorithm in (83) takes the following simple form

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - (\Phi \mathbf{D} \Phi^T)^{-1} \Phi (\mathbf{Q} - \mathbf{Y}) \quad (84)$$

$$= (\Phi \mathbf{D} \Phi^T)^{-1} (\Phi \mathbf{D} \Phi^T \mathbf{w}_t - \Phi (\mathbf{Q} - \mathbf{Y})) \quad (85)$$

$$= (\Phi \mathbf{D} \Phi^T)^{-1} \Phi \mathbf{D} \mathbf{z}, \quad (86)$$

where  $\mathbf{D} \in \mathbb{R}^{N \times N}$  is the diagonal matrix defined in Exercise 2.6.1 and we have denoted  $\mathbf{z} := \Phi^T \mathbf{w}_t - \mathbf{D}^{-1} (\mathbf{Q} - \mathbf{Y}) \in \mathbb{R}^N$ . From the analysis in Exercise 1.2.1, one can observe that the last expression is the solution of the following ‘weighted least squares problem’

$$\tilde{\mathbf{w}} := \arg \min_{\mathbf{w} \in \mathbb{R}^p} \|\mathbf{z} - \sqrt{\mathbf{D}} \Phi^T \mathbf{w}\|_F^2. \quad (87)$$

However, both the target vector  $\mathbf{z}$  and the weight matrix  $\sqrt{\mathbf{D}}$  change when we update  $\mathbf{w}_t$ . Hence one has to apply the above formula iteratively. For these reasons, the Newton-Ralphson algorithm (83) for logistic regression is also called *iterative reweighted least squares* (IRWS).

**Exercise 2.2.2.** The gradient descent (GD) algorithm for logistic regression training in (82) is implemented as the python function ‘fit\_LR\_GD’ in the [Jupyter notebook](#) provided in the course repository.

```
def fit_LR_GD(Y, H, W0=None, sub_iter=100, stopping_diff=0.01):
    """
    Convex optimization algorithm for Logistic Regression using Gradient Descent
    Y = (n x 1), H = (p x n) (\Phi in lecture note), W = (p x 1)
    Logistic Regression: Y ~ Bernoulli(Q), Q = sigmoid(H.T @ W)
    MLE -->
    Find \hat{W} = argmin_W ( sum_j ( log(1+exp(H_j.T @ W)) - Y.T @ H.T @ W ) )
    ...
    if W0 is None:
        W0 = np.random.rand(H.shape[0],1) #If initial coefficients W0 is None, randomly initialize
    W1 = W0.copy()
    i = 0
    dist = 1
    grad = np.ones(W0.shape)
    while (i < sub_iter) and (np.linalg.norm(grad) > stopping_diff):
        Q = 1/(1+np.exp(-H.T @ W1)) # probability matrix, same shape as Y
        # grad = H @ (Q - Y).T + alpha * np.ones(W0.shape[1])
        grad = H @ (Q - Y)
        W1 = W1 - (np.log(i+1) / (((i + 1) ** (0.5)))) * grad
        i = i + 1
        # print('iter %i, grad_norm %f' %(i, np.linalg.norm(grad)))
    return W1
```

FIGURE 12. A python implementation of the gradient descent algorithm for logistic regression training (82).

<sup>1</sup>See Figure 1: The trajectory on the graph of the loss function can be viewed as that of a gradient descent algorithm.

<sup>2</sup>This is an instance of dinimishing step sizes, where the step sizes are required to be non-summable  $\sum_{t=1}^{\infty} \eta_t = \infty$  but square-summable  $\sum_{t=1}^{\infty} \eta_t^2 < \infty$ . See [SSY18].

<sup>3</sup>The logic behind this algorithm is to minimize the second-order Taylor approximation of the objective  $\ell_{LR}$ . Namely,  $\ell_{LR}(\mathbf{w}) = \ell_{LR}(\mathbf{w}_t) + \nabla \ell_{LR}(\mathbf{w}_t)^T (\mathbf{w} - \mathbf{w}_t) + \frac{1}{2} (\mathbf{w} - \mathbf{w}_t)^T \nabla^2 \ell_{LR}(\mathbf{w}_t) (\mathbf{w} - \mathbf{w}_t) + O(\|\mathbf{w} - \mathbf{w}_t\|_F^3)$ , and the minimizer of the quadratic function here is exactly (83).

- (i) Modify the code (or make your own function) so that each gradient descent step (82), the current value of objective function  $\ell_{LR}$  in (78) is printed. Test it on the MNIST example with digits [0', 1']. Does the loss monotonically decrease?
- (ii) The implemented version of GD for (82) uses the step size  $\eta_t = \gamma \log(t+1)/\sqrt{(t+1)}$  with  $\gamma = 1$ . Test the same step sizes for  $\gamma = 10$  and also  $\eta_t = 1$  and  $\eta_t = 1/t$ . Test it on the MNIST example with digits [0', 1']. Is there any difference in the way the loss function decrease?
- (iii) If we use the step size of the form  $\eta_t = \gamma t^{-\delta}$  for  $0 \leq \delta \leq 1$ , what would be your optimal choice of the hyperparameters  $\gamma$  and  $\delta$ ? (There is no answer and trying out a few choices and make your observation.)

**Exercise 2.2.3.** See the [Jupyter notebook](#) provided in the course repository.

```
def fit_LR_NR(Y, H, W0=None, sub_iter=100, stopping_diff=0.01):
    """
    Convex optimization algorithm for Logistic Regression using Newton-Ralphson algorithm.
    Y = (n x 1), H = (p x n) (\Phi in lecture note), W = (p x 1)
    Logistic Regression: Y ~ Bernoulli(Q), Q = sigmoid(H.T @ W)
    MLE -->
    Find \hat{W} = argmin_W ( sum_j ( log(1+exp(H_j.T @ W)) - Y.T @ H.T @ W ) )
    """
    ### Implement by yourself.
```

FIGURE 13. A python implementation of the Newton-Ralphson algorithm for logistic regression training (83) to be completed.

- (i) Implement the Newton-Ralphson algorithm for logistic regression training in (83) in the same notebook, and reproduce Figure 23. How does the result compare with two optimization algorithms?
- (ii) Compare with the gradient descent algorithm. Which one do you think is faster? Which one do you think is more efficient? (Do some experiments make your own argument. There is no single answer.)

### 3. Binary classification of MNIST by Logistic Regression

**3.1. Training.** Recall that the MNIST dataset (see Figure 10) consists of total 70000 black-and-white images of size  $28 \times 28$  pixels, corresponding to one of the 10 digits from  $\{0, 1, \dots, 9\}$ . In this subsection, we will apply logistic regression to classify the MNIST images of chosen two digits.

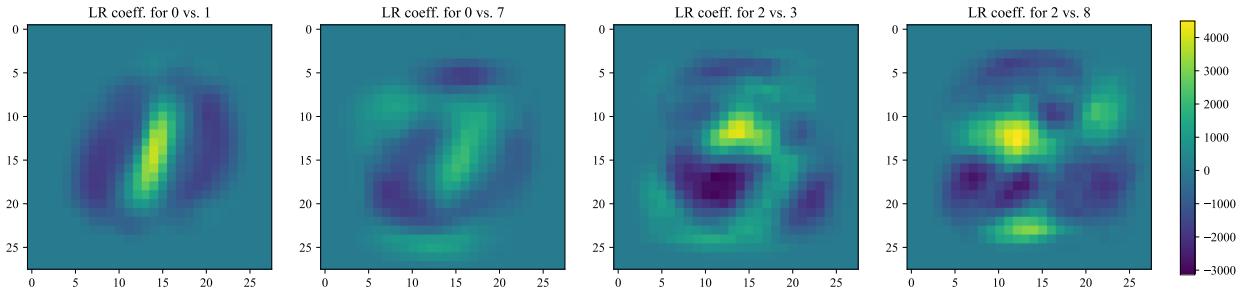


FIGURE 14. Examples of the fitted logistic regression coefficients to the [MNIST](#) dataset with chosen two digits. Gradient descent algorithm is used for training.

For instance, suppose we choose the images for digits '0' and '1'. There are total 6903 and 7877 images of digits '0' and '1', respectively. Viewing each image as a matrix in  $\mathbb{R}^{28 \times 28}$ , the binary classification problem now is to learn a decision algorithm  $\mathcal{A} : \mathbb{R}^{28 \times 28} \supseteq \{\text{Images of '0' or '1'}\} \rightarrow \{0, 1\}$ , where class label 1 means the image is classified as digit '1' and label 0 for digit '0'. The goal is to learn the optimal logistic regression parameter  $\hat{w} \in \mathbb{R}^{28 \times 28}$  by solving the optimization problem (78). The meaning of the  $(i, j)$  entry  $\hat{w}[i, j]$  is the positive association of the  $(i, j)$  pixel of all  $28 \times 28$  images in the train set with class label 1.

The first subplot in Figure 23 displays  $\hat{w}$  for digits '0' and '1', learned by the gradient descent algorithm in (82) with  $\gamma = 1$ . One can see that the pixels lying on the overall shape of digit '1' have large regression coefficients, meaning that they contribute to the image being classified as '1' with a large probability. On the other hand, one can also see that the pixels lying on the overall shape of digit '0' have small regression coefficients, meaning that they discount the probability of the image being classified as '1' and hence encourage it to be classified as '0'. The results for other digit pairs ['0', '7'], ['2', '3'], and ['2', '8'] are also shown in Figure 23.

Notice such a pixel-based regression for image classification would make sense only if the images are well-aligned. For example, if the images for digit '1' show up in all different locations and orientations and even with various scaling, then it won't be possible to attribute any pixel (location in the image) as an important feature of being class 1.

**Exercise 2.3.1.** See the [Jupyter notebook](#) provided in the course repository.

- Reproduce the experiments in Figure 23 for all pairs of digits. (There will be  $\binom{10}{2} = 45$  pairs. Use either convex optimization algorithms.)
- By looking at the plot of the optimal regression coefficients, for each digit  $i$ , make a guess on digits  $j$  that would be the easiest and the hardest for the logistic regression model to distinguish.

**3.2. Testing.** Now let's discuss testing our trained logistic regression model for the MNIST dataset. We first need to discuss how we measure the accuracy of binary classification. Suppose we have target and predicted binary vectors  $y_{\text{test}}, y_{\text{pred}} \in \{0, 1\}^N$ . For each test example  $i$ , there are four possibilities:  $(y_{\text{test}}[i], y_{\text{pred}}[i]) \in \{(T, T), (T, F), (F, T), (F, F)\}$ , where  $T$  stands for 'true' or 'positive' and  $F$  stands for 'false' or 'negative'. These four cases have the following specific names:

- True Positive*:  $(y_{\text{test}}[i], y_{\text{pred}}[i]) = (T, T)$
- False Positive* ( $y_{\text{test}}[i], y_{\text{pred}}[i] = (T, F)$ ) (e.g., 'patient diagnosed to have cancer but does not')
- False Negative* ( $y_{\text{test}}[i], y_{\text{pred}}[i] = (F, T)$ ) (e.g., 'patient diagnosed healthy but has cancer')
- True Negative* ( $y_{\text{test}}[i], y_{\text{pred}}[i] = (F, F)$ )

One of the above instances hold for each test example  $i$ . The table in Figure 15 summarizes various statistics associated to predicting  $y_{\text{test}}$  as  $y_{\text{pred}}$ .

		CONDITION determined by "Gold Standard"		PREVALENCE $\frac{\text{CONDITION POS}}{\text{TOTAL POPULATION}}$	
TOTAL POPULATION		CONDITION POS	CONDITION NEG		
TEST OUT-COME	TEST POS	True Pos TP	Type I Error False Pos FP	Precision Pos Predictive Value $\text{PPV} = \frac{\text{TP}}{\text{TEST P}}$	False Discovery Rate $\text{FDR} = \frac{\text{FP}}{\text{TEST P}}$
	TEST NEG	Type II Error False Neg FN	True Neg TN	False Omission Rate $\text{FOR} = \frac{\text{FN}}{\text{TEST N}}$	Neg Predictive Value $\text{NPV} = \frac{\text{TN}}{\text{TEST N}}$
ACCURACY ACC $\text{ACC} = \frac{\text{TP} + \text{TN}}{\text{TOT POP}}$	Sensitivity (SN), Recall Total Pos Rate $\text{TPR} = \frac{\text{TP}}{\text{CONDITION POS}}$	Fall-Out False Pos Rate $\text{FPR} = \frac{\text{FP}}{\text{CONDITION NEG}}$	Pos Likelihood Ratio $\text{LR}+ = \frac{\text{TPR}}{\text{FPR}}$	Diagnostic Odds Ratio $\text{DOR}$	
	Miss Rate False Neg Rate $\text{FNR} = \frac{\text{FN}}{\text{CONDITION POS}}$	Specificity (SPC) True Neg Rate $\text{TNR} = \frac{\text{TN}}{\text{CONDITION NEG}}$	Neg Likelihood Ratio $\text{LR}- = \frac{\text{TNR}}{\text{FNR}}$	$\text{DOR} = \frac{\text{LR}+}{\text{LR}-}$	

FIGURE 15. Various statistics of binary classification testing. Figure excerpted from [link](#).

Here is a list of a few commonly used statistics:

- $Accuracy = \frac{\# \text{ predicted correctly}}{\text{size of the test set}}$  (e.g. =  $\frac{\# \text{ predicted correctly as cancer or healthy}}{\# \text{ total subjects}}$ )
- $Sensitivity = \frac{\# \text{ predicted positive correctly}}{\# \text{ positive test examples}}$  (e.g. =  $\frac{\# \text{ predicted correctly as cancer}}{\# \text{ cancer patients}}$ )
- $Specificity = \frac{\# \text{ predicted negative correctly}}{\# \text{ negative test examples}}$  (e.g. =  $\frac{\# \text{ predicted correctly as healthy}}{\# \text{ healthy subjects}}$ )
- $Precision = \frac{\# \text{ predicted positive correctly}}{\# \text{ predicted positive}}$  (e.g. =  $\frac{\# \text{ predicted correctly as cancer}}{\# \text{ predicted as cancer}}$ )

Next, before we proceed to the running example of LR on MNIST, we need to say a bit more on how we get the binary vector of predictions  $y_{\text{pred}}$ . Namely, recall that LR is a probabilistic model, where one models the output as a Bernoulli RV  $Y$  with success probability  $p = \sigma(\phi(\mathbf{x})^T \hat{\mathbf{w}})$  as in (71). The result of fitting LR to the training set is the optimal parameter  $\hat{\mathbf{w}}$  defined in (78), so we can compute the predictive probability  $p = p(y|\phi(\mathbf{x}))$  for the unknown output. Suppose we have  $p = 0.8$ . This means that our prediction on the output has probability 0.8 to be 1 and 0.2 to be 0, so it would be reasonable to predict the output as 1. More precisely, for the current example, we use the following decision algorithm

$$\mathbf{x} \mapsto \hat{\mathbf{y}} = \mathcal{A}(\mathbf{x}) = \begin{cases} 1 & \text{if } p = \sigma(\phi(\mathbf{x})^T \hat{\mathbf{w}}) \geq 0.5 \\ 0 & \text{if } p = \sigma(\phi(\mathbf{x})^T \hat{\mathbf{w}}) < 0.5. \end{cases} \quad (88)$$

Simply speaking, we are predicting  $\hat{\mathbf{y}}$  to be more likely class label predicted by the predicted probability, with the classification threshold being 0.5<sup>4</sup>. This is an instance of the more general decision algorithm discussed in Exercise 1.3.2.

MNIST Binary Classification by LR for 4 vs. 7

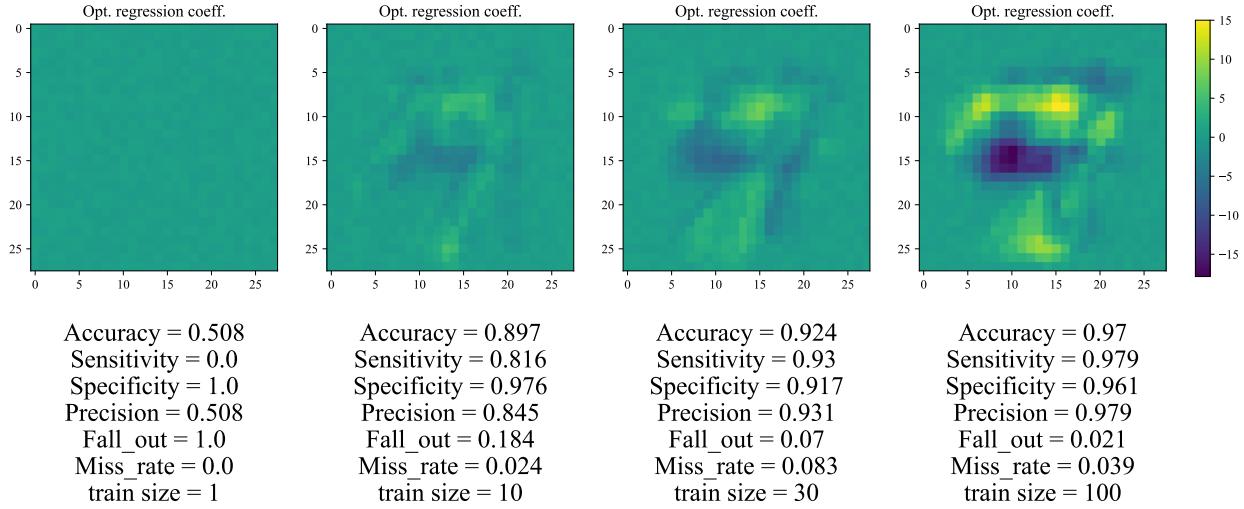


FIGURE 16. Examples of the training and testing logistic regression coefficients to the MNIST dataset with chosen two digits. Gradient descent algorithm is used for training. Training datasets of size in  $\{1, 10, 30, 100\}$  are subsampled and used for fitting the model. Various statistics for binary classification are also shown.

Now we discuss testing LR on the MNIST dataset. MNIST is such a clean dataset so that we only need a few training examples to achieve very good classification accuracies. We test our LR classifier trained on randomly subsampled test sets of sizes 1, 10, 30, and 100 for classifying digits '4' vs. '7'. The test set contains about 20% of examples in the dataset, which is about 14,000 images of digits '4' and '7'. In Figure

<sup>4</sup>The use of threshold 0.5 is not necessary. For highly imbalanced training datasets, the predicted probabilities could be skewed either to 0 or 1, in which case we would need to use some other threshold  $\theta \in [0, 1]$ . The 'overall accuracy' of binary classification using arbitrary threshold  $\theta$  is measured by other means, such as the *area-under-curve* (AUC) score for the corresponding *receiver operation characteristic* (ROC).

[16](#), we have examples of trained regression coefficients (shown as a heat map on  $28 \times 28$  image) as well as various statistics of binary classification we discussed before. Observe that as we increase the size of the train set, the optimal regression coefficients become more aware of the shapes of ‘7’ as positive (yellow) and ‘4’ as negative (blue), and the classification metrics improve in general. See Figure [17](#) for a similar example with digits ‘3’ and ‘8’ instead, which seems to be a bit more difficult case than distinguishing between ‘4’ and ‘7’.

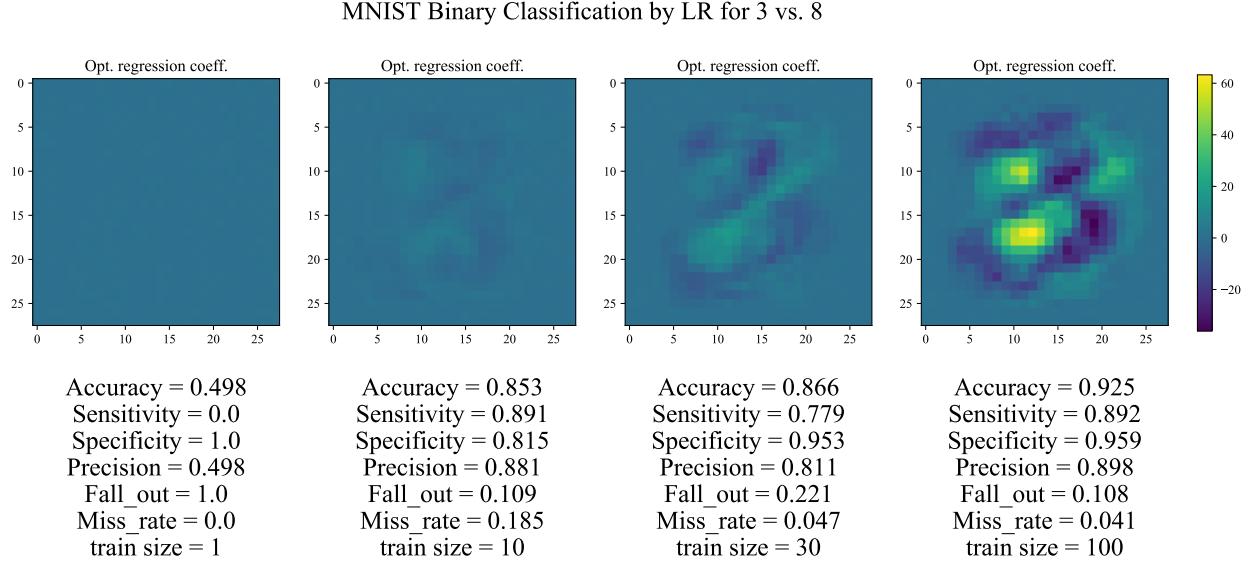


FIGURE 17. Examples of the training and testing logistic regression coefficients to the [MNIST](#) dataset with chosen two digits. Gradient descent algorithm is used for training. Training datasets of size in  $\{1, 10, 30, 100\}$  are subsampled and used for fitting the model. Various statistics for binary classification are also shown.

**3.3. Multiclass Logistic Regression.** Next, we extend our logistic regression model to the multiclass classification, where the class labels now assume  $\kappa \geq 2$  distinct values. This model is called the *multiclass logistic regression* (MLR). Without loss of generality, we can assume that the  $\kappa$  classes are the integers in  $\{1, 2, \dots, \kappa\}$ . Say we have training examples  $(\phi(\mathbf{x}_1), y_1), \dots, (\phi(\mathbf{x}_N), y_N)$ , where

- $\mathbf{x}_1, \dots, \mathbf{x}_N$ : Input data (e.g., collection of all medical records of each patient)
- $\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_N) \in \mathbb{R}^p$ : Features (e.g., some useful (derived) information for each patient)
- $y_1, \dots, y_N \in \{1, 2, \dots, \kappa\}$ :  $\kappa$  class labels (e.g., digits from 0 to 9).

The basic idea of multiclass logistic regression is to model the output  $y$  as a discrete random variable  $Y$  with probability mass function (PMF)  $p = [p_1, \dots, p_\kappa]$  that depends on the observed feature  $\phi(\mathbf{x})$ . More specifically,

$$\mathbb{P}(Y = i | \phi(\mathbf{x})) = p_i, \quad p_i \propto \exp(\phi(\mathbf{x})^T \mathbf{w}_i) \quad \text{for } i = 1, \dots, \kappa. \quad (89)$$

More precisely, the predictive probability distribution  $p = [p_1, \dots, p_\kappa]$  can be written as

$$p_i = \frac{\exp(\phi(\mathbf{x})^T \mathbf{w}_i)}{\sum_{j=1}^{\kappa} \exp(\phi(\mathbf{x})^T \mathbf{w}_j)}, \quad \text{for } i = 1, \dots, \kappa. \quad (90)$$

A probability distribution of the form  $p_i \propto \exp(a_i)$  is called the *softmax distribution* with *activation*  $\mathbf{a} = [a_1, \dots, a_\kappa]$ . Hence in the case of multinomial LR, the predictive probability distribution  $p = [p_1, \dots, p_\kappa]$

is the softmax distribution with ‘linear activation’  $\boldsymbol{\phi}(\mathbf{x})^T \mathbf{w}_1, \dots, \boldsymbol{\phi}(\mathbf{x})^T \mathbf{w}_\kappa$  of the feature vector  $\boldsymbol{\phi}(\mathbf{x})$ . Notice that this model has parameter vectors  $\mathbf{w}_1, \dots, \mathbf{w}_\kappa$ , one for each of the  $\kappa$  class labels<sup>5</sup>.

Next, let’s discuss how to train MLR. We will use the following matrix notations

$$\mathbf{Y} = \begin{bmatrix} \mathbf{1}(y_1 = 1) & \cdots & \mathbf{1}(y_1 = \kappa) \\ \vdots & & \vdots \\ \mathbf{1}(y_N = 1) & \cdots & \mathbf{1}(y_N = \kappa) \end{bmatrix} \in \{0, 1\}^{N \times \kappa}, \quad \boldsymbol{\Phi} = \begin{bmatrix} \uparrow & & \uparrow \\ \boldsymbol{\phi}(\mathbf{x}_1) & \cdots & \boldsymbol{\phi}(\mathbf{x}_N) \\ \downarrow & & \downarrow \end{bmatrix} \in \mathbb{R}^{p \times N}, \quad (92)$$

$$\mathbf{W} = \begin{bmatrix} \uparrow & & \uparrow \\ \mathbf{w}_1 & \cdots & \mathbf{w}_\kappa \\ \downarrow & & \downarrow \end{bmatrix} \in \mathbb{R}^{p \times \kappa}. \quad (93)$$

Again since it is a probabilistic model, we use the method of maximum likelihood. For that, we compute the joint likelihood of observing the output values  $(y_1, \dots, y_N)$  from the tuple of independent RVs  $(Y_1, \dots, Y_N)$  with PMFs  $p_{ij} \propto \exp(\boldsymbol{\phi}(\mathbf{x}_j)^T \mathbf{w}_j)$  where  $p_{ij} = \mathbb{P}(Y_i = j)$ :

$$L(y_1, \dots, y_N; \mathbf{W}) = \mathbb{P}(Y_1 = y_1, \dots, Y_N = y_N; \mathbf{W}) = \prod_{i=1}^N \prod_{j=1}^\kappa (p_{ij})^{\mathbf{1}(y_i=j)} \quad (94)$$

Hence the optimial parameter matrix  $\hat{\mathbf{W}} = [\hat{\mathbf{w}}_1, \dots, \hat{\mathbf{w}}_\kappa] \in \mathbb{R}^{p \times \kappa}$  by

$$\hat{\mathbf{W}} = \arg \max_{\mathbf{W} \in \mathbb{R}^{p \times \kappa}} \log L(y_1, \dots, y_N; \mathbf{W}) \quad (95)$$

$$= \arg \max_{\mathbf{W} \in \mathbb{R}^{p \times \kappa}} \sum_{i=1}^N \sum_{j=1}^\kappa \mathbf{1}(y_i = j) \log p_{ij} \quad (96)$$

$$= \arg \min_{\mathbf{W} \in \mathbb{R}^{p \times \kappa}} \sum_{i=1}^N \sum_{j=1}^\kappa \mathbf{1}(y_i = j) \left[ -\boldsymbol{\phi}(\mathbf{x}_i)^T \mathbf{w}_j + \log \left( \sum_{q=1}^\kappa \exp(\boldsymbol{\phi}(\mathbf{x}_i)^T \mathbf{w}_q) \right) \right] \quad (97)$$

$$= \arg \min_{\mathbf{W} \in \mathbb{R}^{p \times \kappa}} \left[ -\sum_{i=1}^N \sum_{j=1}^\kappa \mathbf{1}(y_i = j) \boldsymbol{\phi}(\mathbf{x}_i)^T \mathbf{w}_j + \sum_{i=1}^N \log \left( \sum_{q=1}^\kappa \exp(\boldsymbol{\phi}(\mathbf{x}_i)^T \mathbf{w}_q) \right) \right], \quad (98)$$

where we have used the following observation

$$1 = \sum_{j=1}^\kappa \mathbf{1}(y_i = j) \quad \text{for all } i = 1, \dots, N. \quad (99)$$

Using a matrix formulation, we have

$$\hat{\mathbf{W}} = \arg \min_{\mathbf{W} \in \mathbb{R}^{p \times \kappa}} \left[ \ell_{MLR}(\mathbf{W}) := \left( \sum_{i=1}^N \log \left( \sum_{q=1}^\kappa \exp(\boldsymbol{\phi}(\mathbf{x}_i)^T \mathbf{w}_q) \right) \right) - \text{tr}(\mathbf{Y}^T \boldsymbol{\Phi}^T \mathbf{W}) \right], \quad (100)$$

The loss function  $\ell_{MLR}(\mathbf{W})$  in (78) is called the *multiclass logistic regression loss*.

In Exercise 2.3.2, we show that the above is a convex optimization problem and also compute the gradient and Hessian. As in the logistic regression case, we can solve Exercise 2.3.2 iteratively either by gradient descent or Newton-Ralphson. For gradient descent algorithm, we use

$$\mathbf{W}_{t+1} \leftarrow \mathbf{W}_t - \eta_t \nabla_{\mathbf{W}} \ell_{MLR}(\mathbf{W}_t) = \mathbf{W}_t - \eta_t \boldsymbol{\Phi}(\mathbf{Q}(\mathbf{W}_t) - \mathbf{Y}) \in \mathbb{R}^{p \times \kappa}, \quad (101)$$

where we can use the same step sizes  $\eta_t = \gamma(\log t)/\sqrt{t}$  for some  $\gamma > 0$ . The Newton-Ralphson algorithm (83) can also be implemented for the case of MLR by using the Hessian computed in Exercise 2.3.2.

---

<sup>5</sup>In fact, it is enough to determine  $\kappa - 1$   $p$ -dimensional vectors  $\mathbf{w}_2 - \mathbf{w}_1, \dots, \mathbf{w}_\kappa - \mathbf{w}_1$  since by dividing both the numerator and the denominator by  $\exp(\boldsymbol{\phi}(\mathbf{x})^T \mathbf{w}_1)$ , we can write the MLR predictive distribution in (90) as

$$p_i = \frac{\exp(\boldsymbol{\phi}(\mathbf{x})^T (\mathbf{w}_i - \mathbf{w}_1))}{\sum_{j=1}^\kappa \exp(\boldsymbol{\phi}(\mathbf{x})^T (\mathbf{w}_j - \mathbf{w}_1))}, \quad \text{for } i = 1, \dots, \kappa \quad (91)$$

**Exercise 2.3.2** (Convexity of the multiclass logistic regression loss). Let  $\ell_{MLR}(\mathbf{W})$  denote the multiclass logistic regression loss defined in (100). Let  $\mathbf{Q} = \mathbf{Q}(\mathbf{W}) = (p_{ij})_{i,j} \in \mathbb{R}^{N \times \kappa}$  where  $p_{ij} \propto \exp(\boldsymbol{\phi}(\mathbf{x}_i)^T \mathbf{w}_j)$ .

(i) Show that

$$\nabla_{\mathbf{W}} \ell_{MLR}(\mathbf{W}) = \boldsymbol{\Phi}(\mathbf{Q} - \mathbf{Y}) \in \mathbb{R}^{p \times \kappa}. \quad (102)$$

(ii) For each  $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_\kappa] \in \mathbb{R}^{p \times \kappa}$ , let  $\mathbf{W}^{\text{vec}} := [\mathbf{w}_1^T, \dots, \mathbf{w}_\kappa^T]^T \in \mathbb{R}^{p\kappa}$  denote its vectorization. Identify  $\mathbf{W}$  with  $\mathbf{W}^{\text{vec}}$  so that  $\ell_{MLR} : \mathbb{R}^{p\kappa} \rightarrow \mathbb{R}$  is a function that maps a vector to a scalar. Then its Hessian  $\mathbf{H} := \frac{\partial^2 \ell_{MLR}(\mathbf{W})}{\partial \mathbf{W} \partial \mathbf{W}^T}$  is a  $p\kappa \times p\kappa$  matrix. Show that

$$\nabla_{\mathbf{W}} \ell_{MLR}(\mathbf{W}) = (\boldsymbol{\Phi}(\mathbf{Q} - \mathbf{Y}))^{\text{vec}} = \underbrace{\begin{bmatrix} \boldsymbol{\Phi} & O & O & \cdots & O \\ O & \boldsymbol{\Phi} & O & \cdots & O \\ \vdots & \vdots & & & \vdots \\ O & \cdots & O & \boldsymbol{\Phi} & O \\ O & O & \cdots & O & \boldsymbol{\Phi} \end{bmatrix}}_{\kappa \times \kappa \text{ block matrix}} (\mathbf{Q}^{\text{vec}} - \mathbf{Y}^{\text{vec}}) \in \mathbb{R}^{p\kappa}, \quad (103)$$

$$\mathbf{H} = \frac{\partial^2 \ell_{MLR}(\mathbf{W})}{\partial \mathbf{W} \partial \mathbf{W}^T} = \nabla_{\mathbf{W}} \mathbf{Q}^T \boldsymbol{\Phi}^T = \begin{bmatrix} \boldsymbol{\Phi} \mathbf{D}_1 \boldsymbol{\Phi}^T & O & \cdots & O \\ O & \boldsymbol{\Phi} \mathbf{D}_2 \boldsymbol{\Phi}^T & \cdots & O \\ \vdots & \vdots & & \vdots \\ O & \cdots & O & \boldsymbol{\Phi} \mathbf{D}_\kappa \boldsymbol{\Phi}^T \end{bmatrix} \in \mathbb{R}^{p\kappa \times p\kappa}, \quad (104)$$

where for  $j = 1, \dots, \kappa$ ,  $\mathbf{D}_j$  is the  $N \times N$  diagonal matrix with diagonal entries  $\mathbf{D}_{jj} = p_{jj}(1 - p_{jj})$ .

(iii) Argue that the Hessian  $\mathbf{H}$  in (ii) is positive semi-definite. Conclude that the logistic regression loss  $\ell_{LR}$  is convex. (Hint: Each  $\mathbf{D}_j$  is a diagonal matrix with positive diagonal entries.)

#### 4. Classification of MNIST by Multiclass Logistic Regression

**4.1. Training.** Recall that the MNIST dataset (see Figure 10) consists of total 70000 black-and-white images of size  $28 \times 28$  pixels, corresponding to one of the 10 digits from  $\{0, 1, \dots, 9\}$ . In this subsection, we will apply multiclass logistic regression to classify the MNIST images of more than two digits.

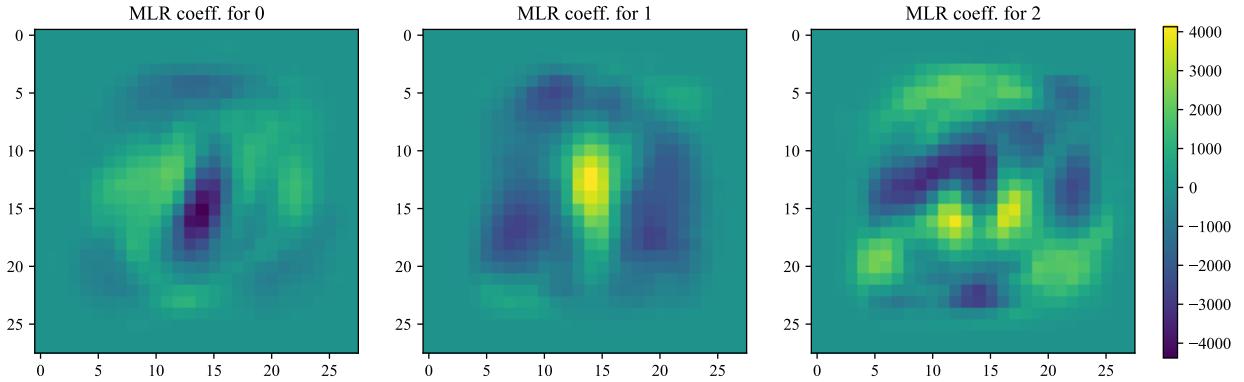


FIGURE 18. Examples of the fitted logistic regression coefficients to the MNIST dataset with digits in  $[0', 1', 2']$ . Gradient descent algorithm is used for training.

For instance, suppose we choose the images for digits '0', '1', and '2'. There are total 6903, 7877, and 6990 images of digits '0', '1', and '2', respectively. Viewing each image as a matrix in  $\mathbb{R}^{28 \times 28}$ , the 3-class classification problem now is to learn a decision algorithm  $\mathcal{A} : \mathbb{R}^{28 \times 28} \supseteq \{\text{Images of '0' or '1' or '2'}\} \rightarrow \{0, 1, 2\}$ , where class label  $i$  means the image is classified as digit ' $i$ ' for  $i \in \{0, 1, 2\}$ . The goal is to learn the optimal logistic regression parameter  $\hat{\mathbf{W}} \in \mathbb{R}^{28 \times 28 \times 3}$  by solving the optimization problem (100). The

meaning of the  $(i, j, k)$  entry  $\hat{\mathbf{W}}[i, j, k]$  is the positive association of the  $(i, j)$  pixel of all  $28 \times 28$  images in the train set with class label  $k$ .

The first subplot in Figure 18 displays  $\hat{\mathbf{W}}[:, :, 0]$  learned by the gradient descent algorithm in (101) with  $\gamma = 1$ . One can see that the pixels lying on the overall shape of digit '0' have large (yellow) regression coefficients, meaning that they contribute to the image being classified as '0' with a large probability. On the other hand, one can also see that the pixels lying on the overall shape of digits '1' or '2' have small regression coefficients, meaning that they discount the probability of the image being classified as '0'. The other subplots show  $\hat{\mathbf{W}}[:, :, 1]$  and  $\hat{\mathbf{W}}[:, :, 2]$ .

**4.2. Testing.** Here we discuss testing our trained multiclass logistic regression model for the MNIST dataset. Suppose we have target and predicted vectors  $y_{\text{test}}, y_{\text{pred}} \in \{1, 2, \dots, \kappa\}^N$ . We measure the overall accuracy of this prediction by the *confusion matrix*, which is a  $\kappa \times \kappa$  matrix whose  $(i, j)$  entry means the number of instances of true class  $i$  classified as class  $j$ . Note that in the ideal situation of perfect prediction  $y_{\text{pred}} = y_{\text{test}}$ , the confusion matrix will be  $N$  times the  $\kappa \times \kappa$  identity matrix.

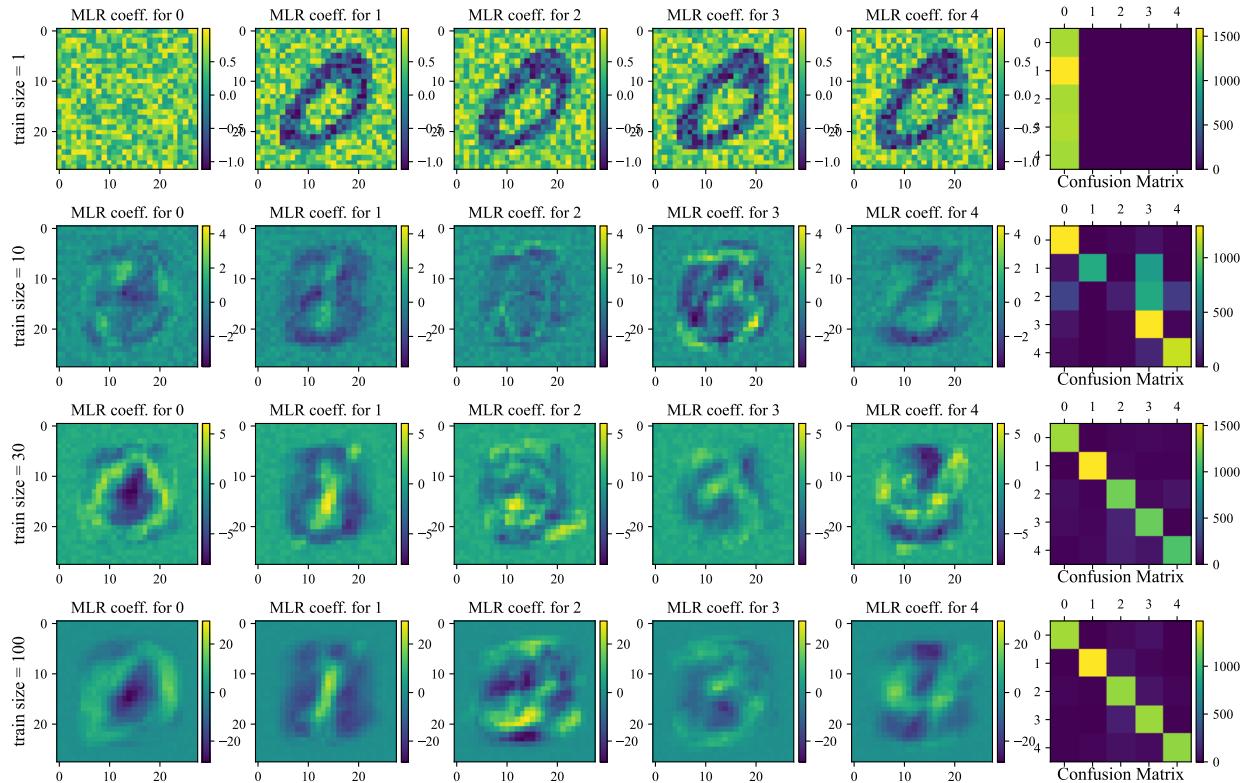


FIGURE 19. Examples of the fitted logistic regression coefficients to the [MNIST](#) dataset with digits in  $[0', 1', 2', 3', 4']$ . Gradient descent algorithm is used for training.

Next, recall that MLR is a probabilistic model, where one models the output as a discrete RV  $Y$  with PMF  $p = [p_1, \dots, p_\kappa]$ ,  $p_i \propto \exp(\boldsymbol{\phi}(\mathbf{x})^T \mathbf{w}_i)$  as in (89). The result of fitting MLR to the training set is the optimal parameter  $\hat{\mathbf{W}} = [\hat{\mathbf{w}}_1, \dots, \hat{\mathbf{w}}_\kappa]$  defined in (100), so we can compute the predictive PMF  $p = p(\mathbf{y}|\boldsymbol{\phi}(\mathbf{x}))$  for the unknown output. We will use the following *maximum a posteriori* decision algorithm

$$\mathbf{x} \mapsto \hat{\mathbf{y}} = \mathcal{A}(\mathbf{x}) = \arg \max_{i \in \{1, 2, \dots, \kappa\}} p(i | \boldsymbol{\phi}(\mathbf{x})) \quad (105)$$

$$= \arg \max_{i \in \{1, 2, \dots, \kappa\}} \exp(\boldsymbol{\phi}(\mathbf{x})^T \hat{\mathbf{w}}_i) \quad (106)$$

$$= \arg \max_{i \in \{1, 2, \dots, \kappa\}} \boldsymbol{\phi}(\mathbf{x})^T \hat{\mathbf{w}}_i. \quad (107)$$

Simply speaking, we are predicting  $\hat{y}$  to be more likely class label predicted by the predicted PMF.

In Figure 19, we test our MLR classifier trained on randomly subsampled test sets of sizes 1, 10, 30, and 100 for classifying digits between '0', '1', ..., '4'. The test set contains about 20% of examples in the dataset, which is about 35,000 images of the five digits. In Figure 19, we have examples of trained regression coefficients (shown as a heat map on  $28 \times 28$  image) as well as the confusion matrices. Observe that as we increase the size of the train set, the optimal regression coefficients become more aware of the shapes of the correct digits, and the confusion matrix becomes closer to a diagonal matrix, meaning the algorithm gets less 'confused'.

## 5. Probit Regression

**5.1. Model Formulation.** In this section, we discuss a probabilistic regression (PR) model for binary classification similar to logistic regression. The setting is identical, but the 'link function' is not the sigmoid function but something called 'probit' (probability + unit) function. Probit regression is widely used in Economics.

As before, we assume that the two classes are 0 and 1. Let  $(\boldsymbol{\phi}(\mathbf{x}_1), y_1), \dots, (\boldsymbol{\phi}(\mathbf{x}_N), y_N)$  be the training examples, where

- $\mathbf{x}_1, \dots, \mathbf{x}_N$ : Input data (e.g., collection of all medical records of each patient)
- $\boldsymbol{\phi}(\mathbf{x}_1), \dots, \boldsymbol{\phi}(\mathbf{x}_N) \in \mathbb{R}^p$ : Features (e.g., some useful (derived) information for each patient)
- $y_1, \dots, y_N \in \{0, 1\}$ : Binary class labels (e.g., 'normal' or 'pneumonia').

The basic idea of probit regression is to model the output  $y$  as a Bernoulli (0-1) random variable  $Y$  that takes value 1 iff the linear activation  $\boldsymbol{\phi}(\mathbf{x})^T \mathbf{w}$  is at least a threshold  $\theta$ , where  $\mathbf{w} \in \mathbb{R}^p$  is a  $p$ -dimensional parameter vector. The point here is to randomize the  $\theta$  as an independent Gaussian RV. More specifically,

$$Y | \boldsymbol{\phi}(\mathbf{x}) \sim \text{Bernoulli}(p), \quad p = \mathbb{P}(\boldsymbol{\phi}(\mathbf{x})^T \mathbf{w} \geq Z) =: \Psi(\boldsymbol{\phi}(\mathbf{x})^T \mathbf{w}), \quad (108)$$

where  $Z \sim N(0, 1)$  is an independent standard Gaussian RV. We can also write the probit model more compactly as

$$Y | \boldsymbol{\phi}(\mathbf{x}) \sim \mathbf{1}(\boldsymbol{\phi}(\mathbf{x})^T \mathbf{w} \geq Z). \quad (109)$$

Here, the function  $\Psi : x \mapsto \mathbb{P}(x \geq Z) = \mathbb{P}(Z \leq x)$  is called the *probit* function, which is exactly the CDF of standard Gaussian distribution. Similarly as sigmoid function does, the probit function also maps the real line onto the interval [0, 1]. As in Figure 20, we see the probit link maps real numbers to the extreme probabilities 0 and 1 more aggressively than the sigmoid link.

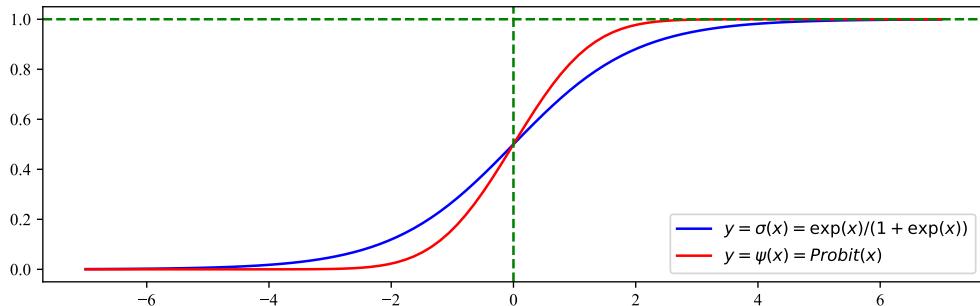


FIGURE 20. Plot of probit and sigmoid function.

Next, let's discuss how to train probit regression model. As before, we will use the method of maximum likelihood. For that, we compute the joint likelihood of observing the output values  $(y_1, \dots, y_N)$

from the tuple of independent Bernoulli variables  $(Y_1, \dots, Y_N)$  with success probabilities  $p_1 := \Psi(\boldsymbol{\phi}(\mathbf{x}_1)^T \mathbf{w}), \dots, p_N := \Psi(\boldsymbol{\phi}(\mathbf{x}_N)^T \mathbf{w})$ :

$$L(y_1, \dots, y_N; \mathbf{w}) = \mathbb{P}(Y_1 = y_1, \dots, Y_N = y_N; \mathbf{w}) = \prod_{i=1}^N p_i^{y_i} (1 - p_i)^{1-y_i}. \quad (110)$$

Hence the optimial parameter  $\hat{\mathbf{w}} \in \mathbb{R}^p$  is given by

$$\hat{\mathbf{w}} = \arg \max_{\mathbf{w} \in \mathbb{R}^p} \log L(y_1, \dots, y_N; \mathbf{w}) \quad (111)$$

$$= \arg \max_{\mathbf{w} \in \mathbb{R}^p} \sum_{i=1}^N y_i \log p_i + (1 - y_i) \log(1 - p_i) \quad (112)$$

$$= \arg \max_{\mathbf{w} \in \mathbb{R}^p} \left[ \sum_{i=1}^N y_i \log \Psi(\boldsymbol{\phi}(\mathbf{x}_i)^T \mathbf{w}) + (1 - y_i) \log(1 - \Psi(\boldsymbol{\phi}(\mathbf{x}_i)^T \mathbf{w})) \right] \quad (113)$$

$$= \arg \min_{\mathbf{w} \in \mathbb{R}^p} \left[ \ell_{PR}(\mathbf{w}) := - \sum_{i=1}^N y_i \log \Psi(\boldsymbol{\phi}(\mathbf{x}_i)^T \mathbf{w}) + (1 - y_i) \log \Psi(-\boldsymbol{\phi}(\mathbf{x}_i)^T \mathbf{w}) \right], \quad (114)$$

where we have use the property  $\Psi(-x) = 1 - \Psi(x)$  of standard normal CDF. The loss function  $\ell_{PR}(\mathbf{w})$  in (114) is called the *probit regression loss*.

In Exercise 2.5.1, we show that the above is a convex optimization problem and also compute the gradient and Hessian. As in the logistic regression case, we can solve Exercise 2.5.1 iteratively either by gradient descent or Newton-Ralphson. For gradient descent algorithm, we use

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta_t \nabla_{\mathbf{w}} \ell_{PR}(\mathbf{w}_t) \in \mathbb{R}^p, \quad (115)$$

where we can use the same step sizes  $\eta_t = \gamma(\log t)/\sqrt{t}$  for some  $\gamma > 0$ . The gradient  $\nabla_{\mathbf{w}} \ell_{PR}(\mathbf{w}_t)$  of the probit loss is computed in Exercise 2.5.1. The Newton-Ralphson algorithm can also be implemented for the case of probit regression similarly by using the Hessian computed in Exercise 2.5.1.

**Exercise 2.5.1** (Convexity of the probit regression loss). Let  $\ell_{PR}(\mathbf{w})$  denote the logistic regression loss defined in (111). Let  $\psi(x) = (2\pi)^{-1/2} \exp(-x^2/2)$  denote the PDF of standard normal distribution and let  $\Psi(x) = \int_{-\infty}^x \psi(t) dt$  denote the corresponding CDF. Let  $\Phi \in \mathbb{R}^{p \times N}$  is the design matrix defined in (79).

(i) Show that

$$\nabla_{\mathbf{w}} \ell_{PR}(\mathbf{w}) = \sum_{i=1}^N \psi(\boldsymbol{\phi}(\mathbf{x}_i)^T \mathbf{w}) \left[ \frac{1 - y_i}{\Psi(-\boldsymbol{\phi}(\mathbf{x}_i)^T \mathbf{w})} - \frac{y_i}{\Psi(\boldsymbol{\phi}(\mathbf{x}_i)^T \mathbf{w})} \right] \boldsymbol{\phi}(\mathbf{x}_i) \quad (116)$$

$$= \Phi(\mathbf{Q}_1 - \mathbf{Q}_2), \quad (117)$$

where

$$\mathbf{Q}_1 = \begin{bmatrix} \frac{(1-y_1)\psi(\boldsymbol{\phi}(\mathbf{x}_1)^T \mathbf{w})}{\Psi(-\boldsymbol{\phi}(\mathbf{x}_1)^T \mathbf{w})} & \dots & \frac{(1-y_N)\psi(\boldsymbol{\phi}(\mathbf{x}_N)^T \mathbf{w})}{\Psi(-\boldsymbol{\phi}(\mathbf{x}_N)^T \mathbf{w})} \end{bmatrix}^T, \quad \mathbf{Q}_2 = \begin{bmatrix} \frac{y_1\psi(\boldsymbol{\phi}(\mathbf{x}_1)^T \mathbf{w})}{\Psi(\boldsymbol{\phi}(\mathbf{x}_1)^T \mathbf{w})} & \dots & \frac{y_N\psi(\boldsymbol{\phi}(\mathbf{x}_N)^T \mathbf{w})}{\Psi(\boldsymbol{\phi}(\mathbf{x}_N)^T \mathbf{w})} \end{bmatrix}^T \in \mathbb{R}^N. \quad (118)$$

(ii) Show the following properties of PDF and CDF of standard normal distribution:

$$\psi(x) = \psi(-x), \quad \Psi(-x) = 1 - \Psi(x), \quad \psi'(x) = -x\psi(x), \quad (119)$$

$$\frac{\partial}{\partial x} \frac{\psi(x)}{\Psi(x)} = \frac{\psi(x)[-x\Psi(x) - \psi(x)]}{\Psi(x)^2}, \quad \frac{\partial}{\partial x} \frac{\psi(x)}{\Psi(-x)} = \frac{\psi(x)[-x\Psi(-x) + \psi(x)]}{\Psi(-x)^2}. \quad (120)$$

(iii) Show that

$$\frac{\partial(\mathbf{Q}_1^T - \mathbf{Q}_2^T)}{\partial \mathbf{w}} = \begin{bmatrix} \uparrow & & \uparrow \\ G_1 \boldsymbol{\phi}(\mathbf{x}_1) & \dots & G_N \boldsymbol{\phi}(\mathbf{x}_N) \\ \downarrow & & \downarrow \end{bmatrix} = \Phi \mathbf{D}, \quad (121)$$

where for  $z_i := \boldsymbol{\phi}(\mathbf{x}_i)^T \mathbf{w}$ , we denote

$$G_i := (1 - y_i) \frac{\psi(z_i) [-z_i \Psi(-z_i) + \psi(z_i)]}{\Psi(-z_i)^2} + y_i \frac{\psi(z_i) [z_i \Psi(z_i) + \psi(z_i)]}{\Psi(z_i)^2}. \quad (122)$$

and  $\mathbf{D}$  is the  $N \times N$  diagonal matrix with diagonal entries  $\mathbf{D}_{ii} = G_i$ .

(iv) Show that

$$\mathbf{H} := \frac{\partial^2 \ell_{PR}(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^T} = \nabla_{\mathbf{w}} (\mathbf{Q}_1^T - \mathbf{Q}_2^T) \boldsymbol{\Phi}^T = \boldsymbol{\Phi} \mathbf{D} \boldsymbol{\Phi}^T \in \mathbb{R}^{p \times p}. \quad (123)$$

(v) (Optional\*) Argue that the Hessian  $\mathbf{H}$  is positive semi-definite. Conclude that the probit regression loss  $\ell_{PR}$  is convex. (Hint:  $\mathbf{D}$  is a diagonal matrix with positive diagonal entries. This was clear for logistic regression. Here, one can show the functions  $-z\Psi(-z) + \psi(z)$  and  $-z\Psi(z) + \psi(z)$  always assume positive values by some calculus argument. See [ZL12, Lem. 1])

## 6. Binary classification of MNIST by Probit Regression

We revisit the MNIST data with probit regression. In this subsection, we will apply probit regression to classify the MNIST images of chosen two digits. The first subplot in Figure 21 displays  $\hat{\mathbf{w}}$  for digits '0' and '1', learned by the gradient descent algorithm in (115) with  $\gamma = 1$ . A similar discussion for LR we made below Figure 23 holds as well in the PR case.

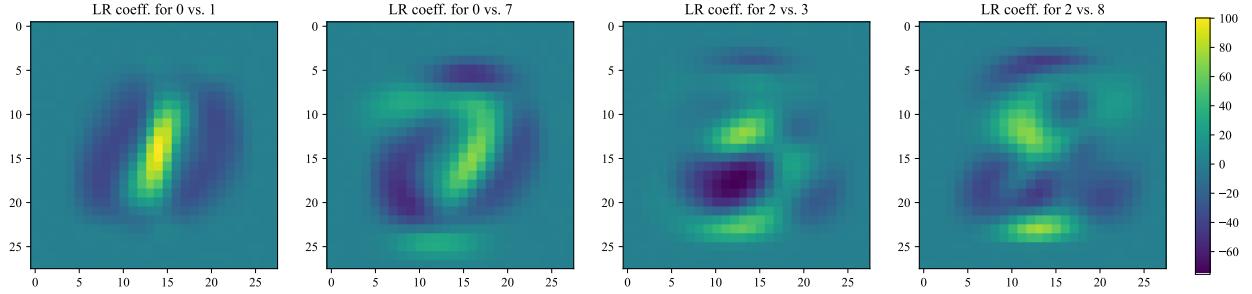


FIGURE 21. Examples of the fitted probit regression coefficients to the MNIST dataset with chosen two digits. Gradient descent algorithm is used for training.

Next, we discuss testing probit regression for binary classification of MNIST dataset. Recall that probit regression is a probabilistic model, where one models the output as a Bernoulli RV  $Y$  with success probability  $p = \Psi(\boldsymbol{\phi}(\mathbf{x})^T \mathbf{w})$  as in (108). The result of fitting PR to the training set is the optimal parameter  $\hat{\mathbf{w}}$  defined in (111), so we can compute the predictive probability  $p = p(\mathbf{y}|\boldsymbol{\phi}(\mathbf{x}))$  for the unknown output. For the current example, we use the following decision algorithm

$$\mathbf{x} \mapsto \hat{\mathbf{y}} = \mathcal{A}(\mathbf{x}) = \begin{cases} 1 & \text{if } p = \Psi(\boldsymbol{\phi}(\mathbf{x})^T \hat{\mathbf{w}}) \geq 0.5 \\ 0 & \text{if } p = \Psi(\boldsymbol{\phi}(\mathbf{x})^T \hat{\mathbf{w}}) < 0.5. \end{cases} \quad (124)$$

Namely, we are predicting  $\hat{\mathbf{y}}$  to be more likely class label predicted by the predicted probability, with the classification threshold being 0.5.

In Figure 16, we show experiments on testing PR on the MNIST dataset similar to LR in Figure 16. As before, we see that as we increase the size of the train set, the optimal regression coefficients become more aware of the shapes of '7' as positive (yellow) and '4' as negative (blue), and the classification metrics improve in general.

Below, we investigate the gradient descent algorithm (115) in a bit more detail as we did for LR in Exercise 2.6.1. We also compare the performance of LR and PR on MNIST dataset with varying degree of noise in the data.

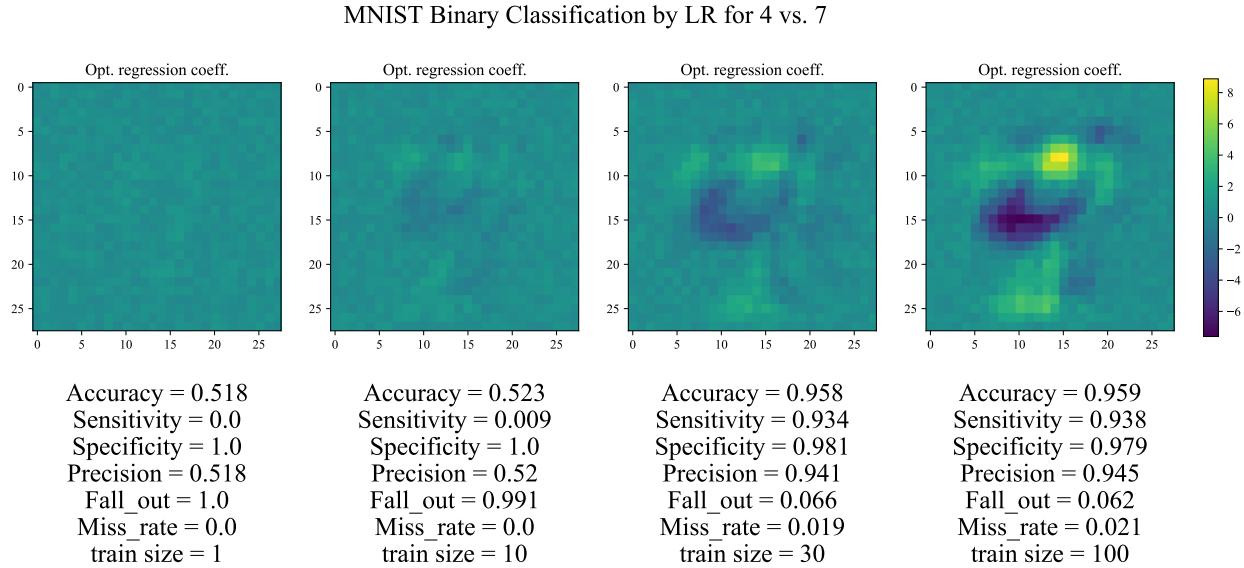


FIGURE 22. Examples of the training and testing probit regression coefficients to the **MNIST** dataset with chosen two digits. Gradient descent algorithm is used for training. Training datasets of size in  $\{1, 10, 30, 100\}$  are subsampled and used for fitting the model. Various statistics for binary classification are also shown.

**Exercise 2.6.1.** The gradient descent (GD) algorithm for probit regression training in (82) is implemented as the python function ‘fit\_PR\_GD’ in the **Jupyter notebook** provided in the course repository.

```
def fit_PR_GD(Y, H, W0=None, sub_iter=100, stopping_diff=0.01):
    """
    Convex optimization algorithm for Probit Regression using Gradient Descent
    Y = (n x 1), H = (p x n) (\Phi in lecture note), W = (p x 1)
    Logistic Regression: Y ~ Bernoulli(Q), Q = Probit(H.T @ W)
    """

    if W0 is None:
        W0 = np.random.rand(H.shape[0], 1) #If initial coefficients W0 is None, randomly initialize

    W1 = W0.copy()
    i = 0
    grad = np.ones(W0.shape)
    while (i < sub_iter) and (np.linalg.norm(grad) > stopping_diff):
        Q = norm.pdf(H.T @ W1) * ((1-Y)/norm.cdf(-H.T @ W1) - Y/norm.cdf(H.T @ W1))
        grad = H @ Q
        W1 = W1 - (np.log(i+1) / (((i + 1) ** (0.5)))) * grad
        i = i + 1
        # print('iter %i, grad_norm %f' %(i, np.linalg.norm(grad)))
    return W1
```

FIGURE 23. A python implementation of the gradient descent algorithm for logistic regression training (115).

- (i) Modify the code (or make your own function) so that each gradient descent step (115), the current value of objective function  $\ell_{PR}$  in (111) is printed. Test it on the MNIST example with digits [‘0’, ‘1’]. Does the loss monotonically decrease?
- (ii) The implemented version of GD for (115) uses the step size  $\eta_t = \gamma \log(t+1)/\sqrt{(t+1)}$  with  $\gamma = 1$ . Test the same step sizes for  $\gamma = 10$  and also  $\eta_t = 1$  and  $\eta_t = 1/(t+1)$ . Test it on the MNIST example with digits [‘0’, ‘1’]. Is there any difference in the way the loss function decrease?
- (iii) If we use the step size of the form  $\eta_t = \gamma t^{-\delta}$  for  $0 \leq \delta \leq 1$ , what would be your optimal choice of the hyperparameters  $\gamma$  and  $\delta$ ? (There is no answer and trying out a few choices and make your observation.)

**Exercise 2.6.2.** In this exercise, we investigate the ‘robustness’ of the logistic and the probit regression for binary classification.

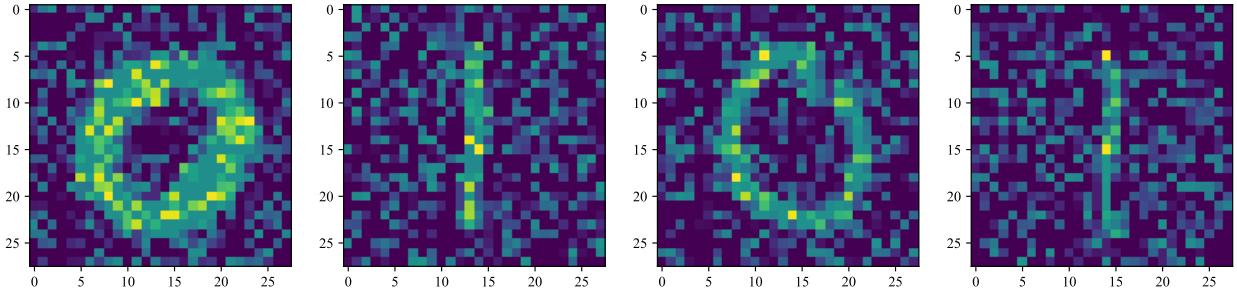


FIGURE 24. Corrupted MNIST images where Uniform([0, 1]) is added to 50% of randomly chosen pixels.

- (i) Modify the data preprocessing function “sample\_binary\_MNIST” so that it has the option to add noise with given ratio. Namely, if ‘noise\_ratio’ is 0.5, then it will add independent Uniform([0, 1]) variables to 50% of randomly chosen pixels. (See Figure 24)
- (ii) For logistic regression, re-generate figures similar to Figure 16 with varying noise ratios in {0.1, 0.5, 0.9}.
- (iii) For probit regression, re-generate figures similar to Figure 22 with varying noise ratios in {0.1, 0.5, 0.9}.
- (iv) From the previous experiments, compare the robustness of LR and PR against random noise you gave. Can you conclude that one is significantly more robust than the other?

## 7. Naive Bayes Classifier

In this section, we discuss a class of classification algorithms called *Naive Bayes classifier*. Recall that in logistic (or probit) regression, we modeled the output class  $Y$  given the feature vector  $\phi(\mathbf{x})$  directly as a RV. In other words, we modeled the *predictive probability*  $p(\text{class}|\text{features})$ . However, Naive Bayes approach, we instead model the likelihood  $p(\text{features}|\text{class})$ , which is also called the *class-conditional probability* in the current situation. This change of perspective is possible by using Bayes’ Theorem, and the ‘Naive Bayes’ assumption is used to simplify the computation of class-conditional probabilities.

**7.1. The model formulation.** We discuss the basic formulation of Naive Bayes Classifier in the setting of multiclass classification. We assume that there are  $\kappa$  class labels in  $\{1, 2, \dots, \kappa\}$ . Say we have training examples  $(\phi(\mathbf{x}_1), y_1), \dots, (\phi(\mathbf{x}_N), y_N)$ , where

- $\mathbf{x}_1, \dots, \mathbf{x}_N$ : Input data (e.g., collection of all medical records of each patient)
- $\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_N) \in \mathbb{R}^p$ : Features (e.g., some useful (derived) information for each patient)
- $y_1, \dots, y_N \in \{1, 2, \dots, \kappa\}$ :  $\kappa$  class labels (e.g., digits from 0 to 9).

As in the multiclass logistic regression, we model the output  $y$  as a discrete random variable  $Y$  with probability mass function (PMF)  $[p_1, \dots, p_\kappa]$  that depends on the observed feature  $\phi(\mathbf{x})$ :

$$\mathbb{P}(Y = i | \phi(\mathbf{X})) = p_i \quad \text{for } i = 1, \dots, \kappa. \quad (125)$$

Instead of modeling the predictive PMF  $p$  directly, we use Bayes’ Theorem to rewrite the above as

$$\underbrace{\mathbb{P}(Y = i | \phi(\mathbf{X}) = [\phi_1, \dots, \phi_m])}_{\text{predictive prob.}} \propto \underbrace{\mathbb{P}(\phi(\mathbf{X}) = [\phi_1, \dots, \phi_p] | Y = i)}_{\text{class-conditional prob.}} \underbrace{\mathbb{P}(Y = i)}_{\text{prior}}. \quad (126)$$

Now the main focus is to model the class-conditional probability. Namely, what is the probability that, given the class label  $Y$  is  $i$ , the observed feature vector  $\phi(\mathbf{X}) = [\phi_1(\mathbf{X}), \dots, \phi_p(\mathbf{X})]$  equals  $[\phi_1, \dots, \phi_p]$ ?

In order to properly model the  $p$ -dimensional random feature vector  $\phi(\mathbf{X}) \in \mathbb{R}^p$ , we need to specify its  $p$ -dimensional class-conditional joint distribution. When the features  $\phi_1, \dots, \phi_p$  assume discrete values,

then the joint distribution can be represented as a  $p \times p$  contingency table, which has  $(p - 1)^2$  degrees of freedom to be specified. When  $p$  is large, such as quadratic dependence of model complexity is not desirable. In order to simplify the situation, we make the following ‘naive Bayes assumption’:

**Naive Bayes Assumption:** Given the class label  $Y = i$ , the features  $\phi(\mathbf{X}) = [\phi_1(\mathbf{X}), \dots, \phi_p(\mathbf{X})]$  are *independent*.

Clearly, such an assumption is not realistic and hence ‘naive’. For instance, the events that a patient has weight  $> 75\text{kg}$  and age  $> 65$  are not independent conditionally on having cancer. However, the naive Bayes assumption simplifies our model to linear complexity:

$$\mathbb{P}([\phi_1(\mathbf{X}), \dots, \phi_p(\mathbf{X})] = [\phi_1, \dots, \phi_p] \mid Y = i) = \prod_{j=1}^p \mathbb{P}(\phi_j(\mathbf{X}) = \phi_j \mid Y = i). \quad (127)$$

Hence, we only need to model the marginal class-conditional probabilities  $\mathbb{P}(\phi_j(\mathbf{X}) = \phi_j \mid Y = i)$  individually and then multiply all to get the joint class-conditional probability.

Suppose we use parameter  $\mathbf{w}_{ij}$  for the marginal class-conditional probability  $\mathbb{P}(\phi_j(\mathbf{X}) = \phi_j(\mathbf{x}_i) \mid Y = y_i)$  of  $j$ th feature given  $i$ th class. Then the joint likelihood of observing the class labels  $y_1, \dots, y_N$  given the observed features  $\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_N)$  can be written as

$$L(y_1, \dots, y_N; \mathbf{W}) = \prod_{s=1}^N \overbrace{\mathbb{P}(Y_s = y_s \mid \phi(\mathbf{X}) = \phi(\mathbf{x}_s))}^{\text{predictive prob.}} \quad (128)$$

$$\propto \prod_{s=1}^N \underbrace{\mathbb{P}(\phi(\mathbf{X}) = \phi(\mathbf{x}_s) \mid Y_s = y_s)}_{\text{class-conditional prob.}} \underbrace{\mathbb{P}(Y_s = y_s)}_{\text{prior}} \quad (129)$$

$$= \prod_{s=1}^N \prod_{i=1}^K \left( \prod_{j=1}^p \mathbb{P}(\phi_j(\mathbf{X}) = \phi_j(\mathbf{x}_s) \mid Y_s = i) \right) \mathbb{P}(Y_s = i)^{\mathbf{1}(y_s=i)}. \quad (130)$$

Hence the optimal parameter  $\hat{\mathbf{W}} = (\hat{\mathbf{w}}_{ij})$  can be found by the standard maximum likelihood method:<sup>6</sup>

$$\hat{\mathbf{W}} = \arg \max_{\mathbf{W}} \left[ \sum_{s=1}^N \sum_{i=1}^K p \mathbf{1}(y_s = i) \log \mathbb{P}(Y_s = i) + \sum_{s=1}^N \sum_{j=1}^p \sum_{i=1}^K \mathbf{1}(y_s = i) \log \underbrace{\mathbb{P}(\phi_j(\mathbf{X}) = \phi_j(\mathbf{x}_s) \mid Y_s = i)}_{\text{parameterized by } \mathbf{w}_{ij}} \right]. \quad (131)$$

In the following sections, we will solve the above optimization problem for some well-known cases of naive Bayes classifiers.

Lastly, given a new input  $\mathbf{x}$  and the optimal parameter  $\hat{\mathbf{W}}$  for naive Bayes classifier found by (131) using the training set, we can make the prediction of the class label using the computed predictive probabilities. As before for the multiclass logistic regression (105), we can use the maximum a posteriori (MAP) decision rule:

$$\mathbf{x} \mapsto \mathcal{A}(\mathbf{x}) = \arg \max_{i \in \{1, 2, \dots, K\}} \overbrace{\mathbb{P}(Y = i \mid \phi(\mathbf{X}) = \phi(\mathbf{x}), \hat{\mathbf{W}})}^{\text{predictive prob.}} \quad (132)$$

$$= \arg \max_{i \in \{1, 2, \dots, K\}} \left[ \sum_{s=1}^N \log \mathbb{P}(Y = i) + \sum_{s=1}^N \sum_{j=1}^p \log \underbrace{\mathbb{P}(\phi_j(\mathbf{X}) = \phi_j(\mathbf{x}_s) \mid Y = i)}_{\text{parameterized by } \hat{\mathbf{w}}_{y_i,j}} \right]. \quad (133)$$

<sup>6</sup>Here we treat prior as given, but we also regard it as a parameter to be optimized.

**7.2. Multinomial Naive Bayes.** In this subsection, we discuss the case naive Bayes classifier with discrete features in the context of document classification. We will make the following version of naive Bayes assumption:

**Naive Bayes Model for Document Classification:**

- A set of ‘vocabulary’  $\mathcal{V} = \{a_1, \dots, a_M\}$  is given. (e.g.,  $\mathcal{V} = \{\text{Dear, Sir, Money, Friend, Thanks}\}$ )
- Class labels  $\mathcal{Y} = \{1, 2, \dots, \kappa\}$  (e.g., 0 = ‘normal’, 1 = ‘spam’)
- For each data point (document)  $\mathbf{x}$ , assume its feature  $\phi(\mathbf{x})$  is given by a frequency vector (histogram)  $\phi(\mathbf{x}) = [c_1, \dots, c_M]$ , where  $c_i$  means the count of word  $a_i$  in  $\mathbf{x}$ .

$$(e.g., \mathbf{x} = \text{"Dear Sir, Money Money Money. Thanks."}, \phi(\mathbf{x}) = [1, 1, 3, 0, 1]) \quad (134)$$

- Given the class label  $Y = i$ , each word in document  $\mathbf{x}$  is sampled independently from the set  $\mathcal{V}$  of vocabulary with probability  $q_{ij}$ .
- Use  $q_i$  for the prior probability of a random document being of class  $i$ .

$$(e.g., \mathbb{P}(\phi(\mathbf{X}) = [1, 1, 3, 0, 1] | Y = \text{spam}) \propto q_{11} \cdot q_{12} \cdot q_{13}^3 \cdot q_{14}^0 \cdot q_{15}). \quad (135)$$

- Let  $\mathbf{W} = (q_i, q_{ij})$  denote the parameter vector.

Suppose we have training examples  $(\phi(\mathbf{x}_1), y_1), \dots, (\phi(\mathbf{x}_N), y_N)$ . For each  $s \in \{1, \dots, N\}$  (data index),  $i \in \{1, \dots, \kappa\}$  (class label), and  $j \in \{1, \dots, M\}$  (word index), using the naive Bayes assumption above, we can write

$$\mathbb{P}(\phi_j(\mathbf{X}) = \phi_j(\mathbf{x}_s) \mid Y = i) = \mathbb{P}(\text{seeing the } j\text{th word } a_j \in \mathcal{V} \text{ } \phi(\mathbf{x}_s) \text{ times in } \mathbf{x}_s \mid Y = i) \quad (136)$$

$$\propto (q_{i,j})^{\phi_j(\mathbf{x}_s)}. \quad (137)$$

Hence we have

$$\hat{\mathbf{W}} = \underset{\mathbf{W}=(q_i, q_{ij})}{\operatorname{argmax}} \left[ \sum_{s=1}^N \sum_{i=1}^{\kappa} M \mathbf{1}(y_s = i) \log q_i + \sum_{s=1}^N \sum_{j=1}^M \sum_{i=1}^{\kappa} \mathbf{1}(y_s = i) \phi_j(\mathbf{x}_s) \log(q_{i,j}) \right] \quad (138)$$

$$= \underset{\mathbf{W}=(q_i, q_{ij})}{\operatorname{argmax}} \left[ M \sum_{i=1}^{\kappa} \underbrace{\left( \sum_{s=1}^N \mathbf{1}(y_s = i) \right)}_{c_i := \# \text{ of class } i \text{ examples}} \log q_i + \sum_{i=1}^{\kappa} \sum_{j=1}^M \underbrace{\left( \sum_{s=1}^N \mathbf{1}(y_s = i) \phi_j(\mathbf{x}_s) \right)}_{c_{i,j} := \# \text{ of } j\text{th word in class } i \text{ examples}} \log(q_{i,j}) \right], \quad (139)$$

where the parameters  $\mathbf{W} = (q_i, q_{ij})$  are subject to the constraint

$$0 \leq q_1, \dots, q_{\kappa} \leq 1, \quad q_1 + \dots + q_{\kappa} = 1 \quad (140)$$

$$0 \leq q_{i,1}, \dots, q_{i,M} \leq 1, \quad q_{i,1} + \dots + q_{i,M} = 1 \quad \text{for } i = 0, 1, \dots, \kappa. \quad (141)$$

In the following proposition, we show that the optimal parameters that solve the above maximum likelihood formulation are given by simple formulas involving empirical frequencies of each class and of each words given each class.

**Proposition 2.7.1.** *The solution  $\hat{\mathbf{W}} = (\hat{q}_i, \hat{q}_{i,j})$  to the optimization problem (138) is given by*

$$\hat{q}_i = \frac{1}{N} \sum_{s=1}^N \mathbf{1}(y_s = i), \quad \hat{q}_{i,j} \propto \sum_{s=1}^N \mathbf{1}(y_s = i) \phi_j(\mathbf{x}_s) \quad \text{for } j \in \{1, \dots, M\}. \quad (142)$$

**PROOF.** First, note that the second term in the objective function in (138) does not depend on the parameters  $q_1, \dots, q_{\kappa}$ . Hence  $[\hat{q}_1, \dots, \hat{q}_{\kappa}]$  is a PMF on  $\{1, \dots, \kappa\}$  that maximizes  $\sum_{i=1}^{\kappa} c_i \log q_i$ , where  $c_i := \sum_{s=1}^N \mathbf{1}(y_s = i) \geq 0$ . Hence the formula for  $\hat{q}_i$ ’s follow from Exercise 2.7.2.

Next, fix a class label  $i \in \{1, \dots, \kappa\}$ . Observe that the only terms in the objective function in (138) that depends on  $q_{i,1}, \dots, q_{i,M}$  is  $\sum_{j=1}^M (\sum_{s=1}^N \mathbf{1}(y_s = i) \phi_j(\mathbf{x}_s)) \log(q_{i,j})$ . Since  $[q_{i,1}, \dots, q_{i,M}]$  is a PMF on  $\{a_1, \dots, a_M\}$ , again the formula for  $\hat{q}_{i,j}$ ’s for  $j = 1, \dots, M$  follow from Exercise 2.7.2.  $\square$

**Exercise 2.7.2.** In this exercise, we solve a simple optimization problem related to the multinomial naive Bayes classifier. Fix a finite set  $\mathcal{X} = \{x_1, \dots, x_M\}$  and real numbers  $c_1, \dots, c_M \geq 0$ . Consider the following optimization problem

$$\hat{\mathbf{w}} = \underset{\mathbf{w}=[w_1, \dots, w_M]^T \in \mathbb{R}^M}{\operatorname{argmax}} \left( L(\mathbf{w}) := \sum_{i=1}^M c_i \log w_i \right) \quad (143)$$

subject to:  $\mathbf{w}$  is a PMF on  $\mathcal{X}$ .

We will show that the solution is given by

$$\hat{\mathbf{w}} = \left[ \frac{c_1}{\sum_{i=1}^M c_i}, \dots, \frac{c_M}{\sum_{i=1}^M c_i} \right]^T. \quad (144)$$

- (i) The constraint set of (143) is the ‘simplex’  $\mathcal{C} = \{\mathbf{w} \in \mathbb{R}^M \mid 0 \leq w_1, \dots, w_M \leq 1, \sum_{i=1}^M w_i = 1\}$ . Denote the larger constraint set  $\bar{\mathcal{C}} = \{\mathbf{w} \in \mathbb{R}^M \mid \sum_{i=1}^M w_i = 1\}$ . Let  $\lambda \in \mathbb{R}$  be a Lagrange multiplier for  $\bar{\mathcal{C}}$ . Then the Lagrangian is

$$g(\lambda, \mathbf{w}) = L(\mathbf{w}) - \lambda \left( \sum_{i=1}^M w_i - 1 \right). \quad (145)$$

Show that  $\lambda$  needs to satisfy

$$\frac{\partial g(\lambda, \mathbf{w})}{\partial \mathbf{w}} = \left[ \frac{c_1}{w_1}, \dots, \frac{c_M}{w_M} \right]^T - \lambda [1, \dots, 1]^T = 0. \quad (146)$$

Conclude that (144) is the global maximum  $L(\mathbf{w})$  over  $\bar{\mathcal{C}}$ .

- (ii) From (i), conclude that (144) is the global maximum  $L(\mathbf{w})$  over  $\mathcal{C}$ .

**Example 2.7.3** (Email Classification). Suppose we have training examples  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{100}, y_{100})$  of normal and spam emails over the set of vocabulary  $\mathcal{V} = \{\text{Dear, Sir, Money, Friend, Thanks}\}$ . The class labels are from  $\{0, 1\}$ , where 0 is for normal and 1 is for spam emails. Let  $c_{i,j}$  denote the total number of  $j$ th word in  $\mathcal{V}$  that appear in all emails in the training set of class  $i$ . Suppose the training set gives the following statistics:

$$c_0 := \sum_{s=1}^{100} \mathbf{1}(y_s = '0') = 70, \quad c_1 := \sum_{s=1}^{100} \mathbf{1}(y_s = '1') = 30, \quad (147)$$

$$[c_{0,1}, c_{0,2}, c_{0,3}, c_{0,4}, c_{0,5}] = [70, 20, 15, 40, 90], \quad (148)$$

$$[c_{1,1}, c_{1,2}, c_{1,3}, c_{1,4}, c_{1,5}] = [73, 65, 98, 20, 120]. \quad (149)$$

Then according to Proposition 2.7.1, the optimal prior and class-conditional probabilities for naive Bayes classifier are given by

$$\text{prior PMF} = [\hat{q}_0, \hat{q}_1] = [0.7, 0.3], \quad (150)$$

$$\text{class-conditional PMF for normal emails} = [\hat{q}_{0,1}, \dots, \hat{q}_{0,5}] = \frac{1}{235} [70, 20, 15, 40, 90], \quad (151)$$

$$\text{class-conditional PMF for spam emails} = [\hat{q}_{1,1}, \dots, \hat{q}_{1,5}] = \frac{1}{376} [73, 65, 98, 20, 120]. \quad (152)$$

Now suppose we have a new email  $\mathbf{x} = \text{“Dear Sir, Money Money Money. Thanks.”}$ , so its bag-of-word feature vector is  $\phi(\mathbf{x}) = [1, 1, 3, 0, 1]$ . In order to classify this email, we compute the predictive probabilities using the trained naive Bayes classifier as follow:

$$\mathbb{P}(Y = 0 \mid \phi(\mathbf{X}) = [1, 1, 3, 0, 1]) \propto \frac{70 \cdot 20 \cdot 15^3 \cdot 40^0 \cdot 90}{235^6} \cdot (0.7) = 1.7674 \times 10^{-6}, \quad (153)$$

$$\mathbb{P}(Y = 1 \mid \phi(\mathbf{X}) = [1, 1, 3, 0, 1]) \propto \frac{73 \cdot 65 \cdot 98^3 \cdot 20^0 \cdot 120}{376^6} \cdot (0.3) = 5.6897 \times 10^{-5}. \quad (154)$$

It follows that

$$\mathbb{P}(Y = 0 | \phi(\mathbf{X}) = [1, 1, 3, 0, 1]) = 0.0301, \quad \mathbb{P}(Y = 1 | \phi(\mathbf{X}) = [1, 1, 3, 0, 1]) = 0.9698. \quad (155)$$

According to the MAP decision rule in (132), we classify  $\mathbf{x}$  as a spam email. ▲

**Exercise 2.7.4.** Consider the following training data: The goal is to predict whether new subjects with

age	income	student	credit_rating	buys_computer
$\leq 30$	high	no	fair	no
$\leq 30$	high	no	excellent	no
[31, 40]	high	no	fair	yes
$> 40$	medium	no	fair	yes
$> 40$	low	no	fair	yes
$> 40$	low	yes	excellent	no
[31, 40]	low	yes	excellent	yes
$\leq 30$	medium	yes	fair	yes
$\leq 30$	low	no	fair	yes
$> 40$	medium	yes	fair	yes
$\leq 30$	medium	yes	excellent	yes
[31, 40]	medium	yes	excellent	yes
[31, 40]	high	no	fair	yes
$> 40$	medium	no	excellent	no

given features (age, income, student, credit rating) would buy a computer, using naive Bayes classifier.

- (i) Compute the maximum likelihood prior distribution on the label  $\{0, 1\}$ , where  $0$  = ‘does not buy computer’ and  $1$  = ‘buys computer’.
- (ii) We will model the feature space as the following 10-dimensional product space of binary values  $\{0, 1\}$ :<sup>7</sup>

$$\phi(\mathbf{x}) = [\mathbf{1}(\text{age} \leq 30), \mathbf{1}(\text{age} \in [31, 40]), \mathbf{1}(\text{age} > 40), \dots, \mathbf{1}(\text{credit} = \text{ecellent})] \in \{0, 1\}^{10}. \quad (156)$$

Compute the maximum likelihood class-conditional probabilities of each feature.

- (iii) Suppose we have the following testing examples:

$$\mathbf{x}_1 = \text{“age } \leq 30, \text{ medium income, student, fair credit rating”} \quad (157)$$

$$\mathbf{x}_2 = \text{“age } \in [31, 40], \text{ low income, not student, fair credit rating”} \quad (158)$$

$$\mathbf{x}_3 = \text{“age } > 40, \text{ high income, not student, excellent credit rating”} \quad (159)$$

Compute the predictive probabilities (posterior distribution of class labels) of each testing examples. Make the prediction. For each example, can you tell which factor affected the most for the classification result?

- (iv) Write a python script (in Jupyter notebook) that implements your computation automatically. Test your code on a similar but larger ( $\geq 100$  training examples and  $\geq 20$  testing examples) dataset of your choice.

**Exercise 2.7.5** (Multinomial Naive Bayes for nonnegative features). In this exercise, we extend the multinomial naive Bayes classifier from discrete features to continuous nonnegative features.

1. Consider a  $\kappa$ -class classification problem for class labels  $\mathcal{Y} = \{1, \dots, \kappa\}$  and the feature vector  $\phi(\mathbf{x}) \in \mathbb{R}^M$  of each data point  $\mathbf{x}$  assume *nonnegative real* values in each coordinate.

---

<sup>7</sup>The resulting naive Bayes classifier is a special instance of the multinomial naive Bayes called *Bernoulli naive Bayes*.

2. We assume the following probabilistic model for (data, class) = ( $\mathbf{X}, Y$ ):

$$(prior) \quad \mathbb{P}(Y = i) =: q_i \quad \text{for } i \in \{1, \dots, \kappa\} \quad (160)$$

$$(\text{class-conditional PMF}) \quad \mathbb{P}(\boldsymbol{\phi}(\mathbf{X}) = [c_1, \dots, c_M] \mid Y = i) \propto \prod_{j=1}^M (q_{i,j})^{c_j} \quad \text{for } i \in \{1, \dots, \kappa\} \quad (161)$$

$$(\text{constraints}) \quad \begin{cases} 0 \leq q_1, \dots, q_\kappa \leq 1, & q_1 + \dots + q_\kappa = 1 \\ 0 \leq q_{i,1}, \dots, q_{i,M} \leq 1, & q_{i,1} + \dots + q_{i,M} = 1 \end{cases} \quad \text{for } i = 0, 1, \dots, \kappa. \quad (162)$$

We denote  $\mathbf{W} = (q_i, q_{i,j})$  for the vector of parameters of our model above.

Now show the following:

- (i) Suppose training examples are given:  $(\boldsymbol{\phi}(\mathbf{x}_s), y_s)$ ,  $s = 1, \dots, N$ , and write  $\boldsymbol{\phi}(\mathbf{x}_s) = [\phi_1(\mathbf{x}_s), \dots, \phi_N(\mathbf{x}_s)]$  for each  $s$ . Show that the maximum likelihood estimate of  $\mathbf{W}$ ,  $\hat{\mathbf{W}}$ , is given by

$$\hat{\mathbf{W}} = \underset{\mathbf{W}=(q_i, q_{i,j})}{\operatorname{argmax}} \left[ M \sum_{i=1}^{\kappa} \left( \sum_{s=1}^N \mathbf{1}(y_s = i) \right) \log q_i + \sum_{i=1}^{\kappa} \sum_{j=1}^M \left( \sum_{s=1}^N \mathbf{1}(y_s = i) \phi_j(\mathbf{x}_s) \right) \log (q_{i,j}) \right]. \quad (163)$$

- (ii) Show that the solution to (163) is given by

$$\hat{q}_i = \frac{1}{N} \sum_{s=1}^N \mathbf{1}(y_s = i), \quad \hat{q}_{i,j} \propto \sum_{s=1}^N \mathbf{1}(y_s = i) \phi_j(\mathbf{x}_s) \quad \text{for } j \in \{1, \dots, M\}. \quad (164)$$

(Hint: Notice that the conclusion of Exercise 2.7.2 holds for nonnegative constants  $c_1, \dots, c_N \geq 0$ , not necessarily only for integers. Then argue similarly as in the proof of Proposition 2.7.1.)

- (iii) Does the same conclusion hold if we allow negative feature values?

## 8. Classification with the 20Newsgroups dataset

All experiments in this section can be found in the Jupyter notebook [Math156\\_20NewsGroups.ipynb](#).

**8.1. Preliminaries.** The 20 Newsgroups data set is a collection of approximately 18,000 newsgroup documents, partitioned (nearly) evenly across 20 different newsgroups. The 20 newsgroups collection has become a popular data set for experiments in text applications of machine learning techniques, such as text classification and text clustering<sup>8</sup>. We will use the convenient [20 Newsgroups loader by sklearn](#) for our experiments.

Of the 20 categories, we will use the following 10 categories:

$$\text{categories} = \left[ \begin{array}{l} \text{'comp.graphics', 'comp.sys.mac.hardware', 'misc.forsale',} \\ \text{'rec.motorcycles', 'rec.sport.baseball', 'sci.med', 'sci.space',} \\ \text{'talk.politics.guns', 'talk.politics.mideast', 'talk.religion.misc'} \end{array} \right] \quad (165)$$

Here is an example of raw text in the 'macs.hardware' category:

```
>>> 4257 th doc: Anyone know what would cause my IIcx to not turn on when I hit
    the keyboard switch? The one in the back of the machine doesn't work either...
    The only way I can turn it on is to unplug the machine for a few minutes, then
    plug it back in and hit the power switch in the back immediately... Sometimes
    this doesn't even work for a long time...
```

I remember hearing about this problem a long time ago, and that a logic  
board failure was mentioned as the source of the problem...is this true?

After we load the raw documents, we need to extract features from them and represent in vector form. There are two widely used vectorization method, the *bag-of-words* (BOW) and the *tf-idf* (term-frequency-inverse-document-frequency). We use a standard vocabulary  $\mathcal{V}$  of 45534 words provided in sklearn. For a given document  $\mathbf{x}$ , its BOW and tf-idf feature vectors are defined by

$$\boldsymbol{\phi}_{\text{BOW}}(\mathbf{x}) = [\# \text{ of } j\text{th word in } \mathbf{x}; 1 \leq j \leq |\mathcal{V}|], \quad (166)$$

---

<sup>8</sup>See <http://qwone.com/~jason/20Newsgroups/>

$$\phi_{\text{tf-idf}}(\mathbf{x}) = \left[ \frac{\# \text{ of } j\text{th word in } \mathbf{x}}{\# \text{ of documents that contains the } j\text{th word}} ; 1 \leq j \leq |\mathcal{V}| \right]. \quad (167)$$

Below are the BOW and tf-idf feature vectors of the sample document (4257th) above.

	Coordinate	Bag-of-words	tf-idf
ago	1286	1	0.115122
back	3503	3	0.288116
board	4850	1	0.138195
cause	6355	1	0.125130
either	12551	1	0.108417
failure	14215	1	0.170555
hearing	17842	1	0.170555
hit	18278	2	0.261832
iicx	19200	1	0.217706
immediately	19354	1	0.155037
keyboard	21813	1	0.160377
know	22116	1	0.075488
logic	23500	1	0.158818
long	23525	2	0.203652
machine	23918	2	0.266555
mentioned	25060	1	0.135408
minutes	25592	1	0.145699
plug	30501	1	0.171330
power	30919	1	0.111667
problem	31391	2	0.203382
remember	33429	1	0.119059
sometimes	37293	1	0.135854
source	37380	1	0.127372
switch	39267	2	0.321834
time	40488	2	0.164800
true	41292	1	0.114800
turn	41456	2	0.269503
unplug	42226	1	0.233554
way	43812	1	0.089654
work	44510	2	0.197874

FIGURE 25. Bag-of-words and tf-idf feature vectors of a sample document from 20Newsgroups.

One can think of the tf-idf representation as an inverse-document-frequency (idf) correction to the BOW representation. Namely, even though some words appear many times in a document (e.g., ‘The’), if it appears in many other documents, we discount its frequency by the idf. On the other hand, if some words appear only a few times in a document (e.g., ‘OECD’), if it does not appear in many other documents, we boost its frequency (relatively to the more widely used words).

**8.2. Classification by multiclass Logistic Regression.** In this subsection, we use the multiclass logistic regression (MLR) in Section 3.3 for the task of document classification for the 20Newsgroups dataset we discussed in Subsection 8.1. In Figure (26), we show the optimal MLR coefficients (for two classes) fitted to the 20Newsgroups dataset using the gradient descent algorithm. We show the coefficients for both choice of vectorizers, BOW and tf-idf. By comparing these plots, we see that the coefficients with BOW vectorizer has larger variance than the one with tf-idf vectorizer (see the reported standard deviation in legends). This is because the inverse-document-frequency acts as a regularizer in feature extraction step and the resulting feature vector has smaller variance (check this numerically).

In Figure 29, we show the confusion matrix of testing the fitted MLR on the 20Newsgroups dataset. As before for the MNIST dataset (see Sec. 4), we use the MAP decision algorithm 105 in order to turn the computed predictive probabilities into final predicted class labels. Overall MLR performs well in the 10-class classification task regardless of the choice of the vectorizer, as the confusion matrices are close to being diagonal.

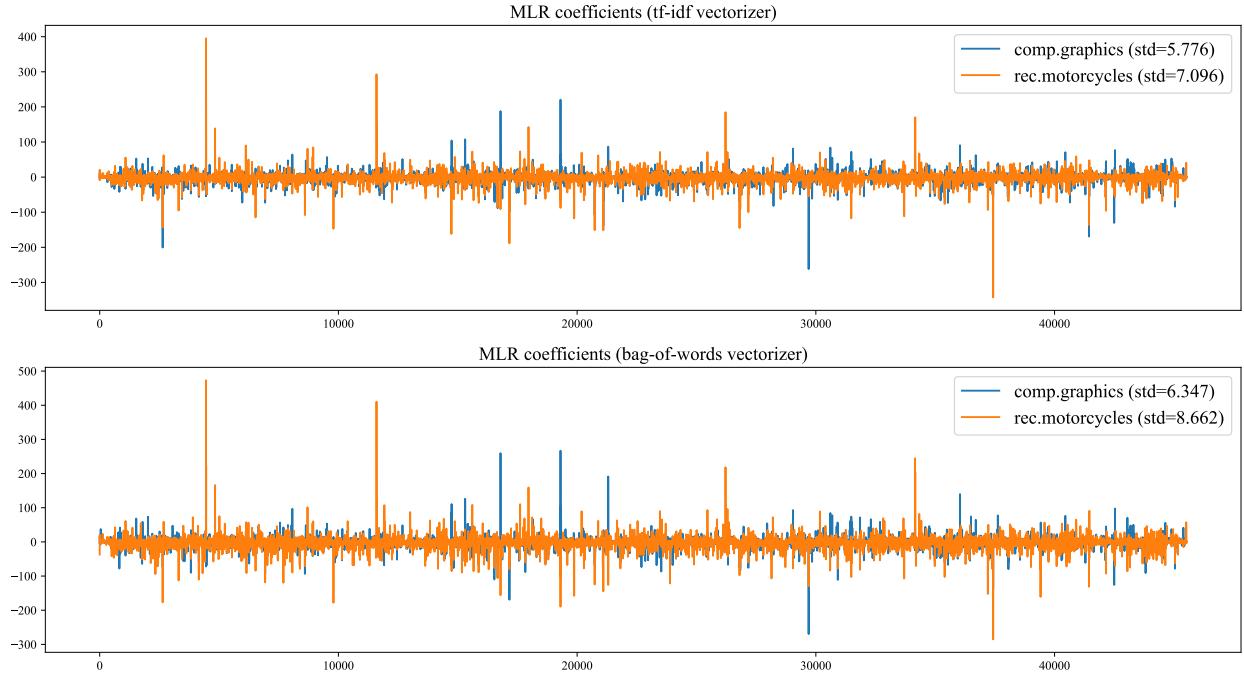


FIGURE 26. Examples of the fitted logistic regression coefficients to the 20Newsgroups dataset with digits. Gradient descent algorithm is used for training. Regression coefficients for the two classes ‘comp.graphics’ and ‘rec.motorcycles’ are shown.

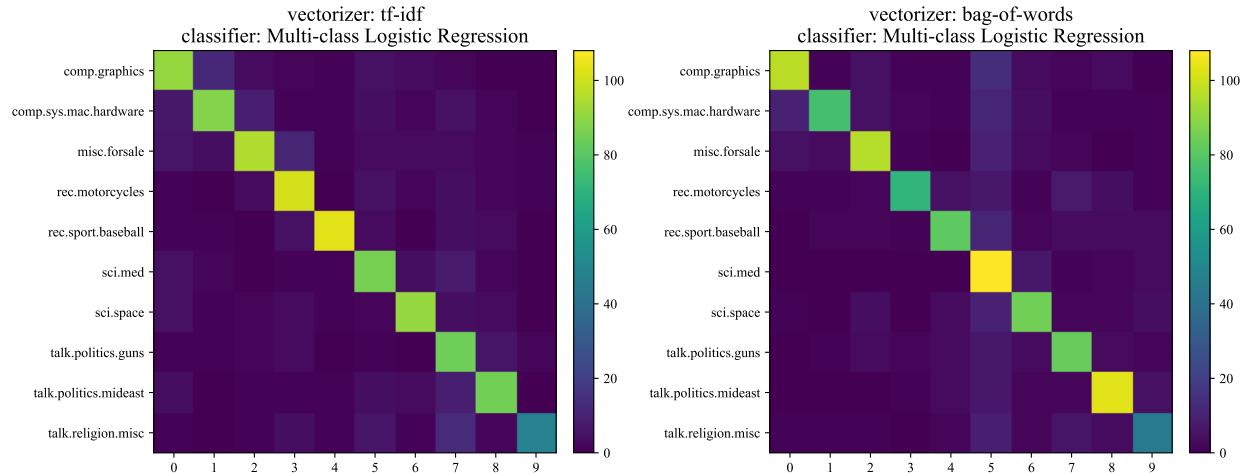


FIGURE 27. Confusion matrices for the 10-class document classification for the 20Newsgroups dataset with multiclass logistic regression. Vectorizers of bag-of-words and tf-idf are both used.

**8.3. Classification by Multinomial Naive Bayes.** In this subsection, we use the multinomial naive Bayesian classifier (MNB) in Section 7.2 for the task of document classification for the 20Newsgroups dataset in Subsection 8.1. In Figure (26), we show the optimal (maximum likelihood) class-conditional probabilities (see Prop. 2.7.1) for two classes. We plot the class-conditional PMFs on the vocabulary for both choice of vectorizers, BOW and tf-idf. Note that even though tf-idf feature vectors are not discrete, according to Exercise 2.7.5, we can still apply MNR to tf-idf feature vectors.

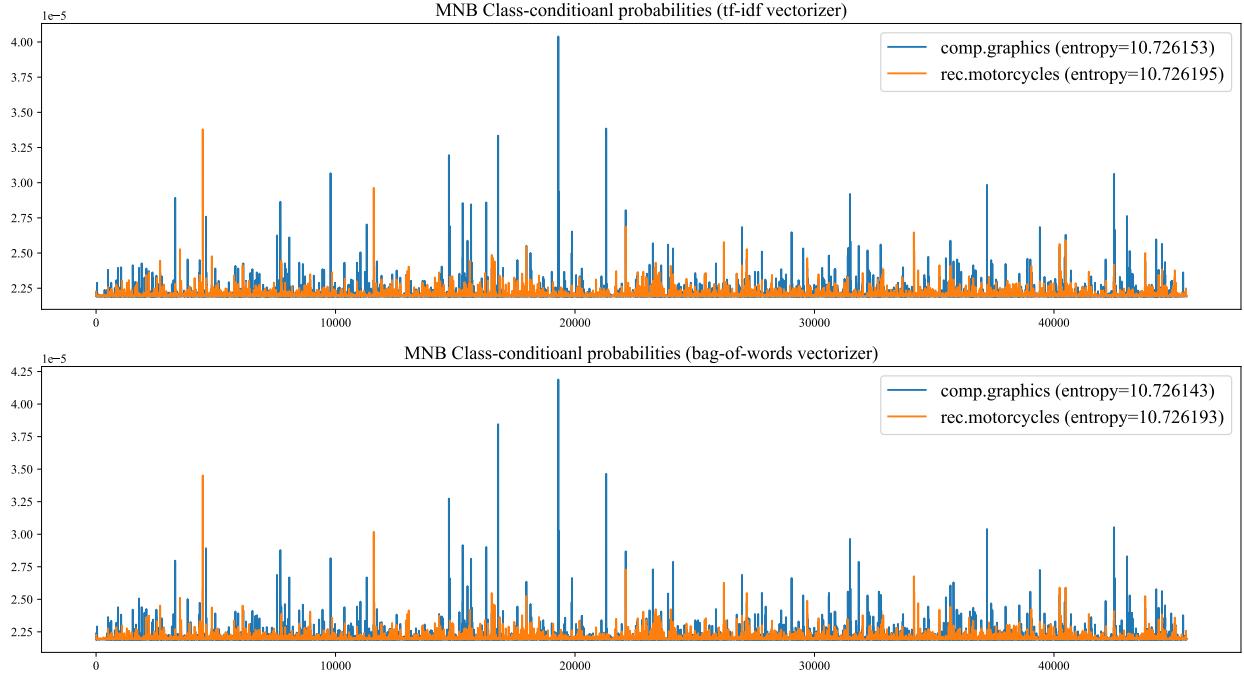


FIGURE 28. Examples of the fitted class-conditional probabilities of the 20Newsgroups dataset with digits. Two classes ‘comp.graphics’ and ‘rec.motorcycles’ are shown.

For each PMF, we also show its (Shannon) *entropy*, which is a measure of the ‘information’ or ‘surprise’ in a probability distribution. For a PMF  $\mathbf{p} = [p_1, \dots, p_n]$ , its entropy  $H(\mathbf{p})$  is defined as

$$H(\mathbf{p}) := - \sum_{i=1}^n p_i \log p_i \geq 0. \quad (168)$$

For any finite sample space, the PMF of maximum entropy on it is the uniform distribution, which should be the ‘most disordered’ among all possible PMFs, so in order to get to know it, we need large samples from it (hence more information content). On the other extreme, the PMFs supported on one single outcome have the lowest possible entropy of zero, and such PMFs contain ‘no surprise’. In general, we regard PMFs with large entropy to be more disordered and hence contain more information. In Figure 28, we see that the class-conditional PMFs for the two classes have almost the same entropy regardless of the choice of the vectorizer. However, the ones with tf-idf vectorizer have slightly larger entropy for both classes. As before, this can be understood as the inverse-document-frequency being a regularizer in feature extraction, so that the feature values are more uniform on the coordinates. Since the (normalised) features have larger entropy, so are the fitted class-conditional probabilities.

In Figure 29, we show the confusion matrix of testing the fitted MNB on the 20Newsgroups dataset. As before for the MNIST dataset (see Sec. 4), we use the MAP decision algorithm 105 in order to turn the computed predictive probabilities into final predicted class labels. Overall MLR performs well in the 10-class classification task regardless of the choice of the vectorizer, as the confusion matrices are close to being diagonal. But less accurately than multiclass logistic regression (see Figure 28). In particular, it seems that MNB suffers classifying documents in the category ‘sci.med’. (Why?)

**Exercise 2.8.1.** According to Exercise 2.7.5, we can apply multinomial naive Bayes for the image classification of MNIST dataset, since the features there (pixel values) are nonnegative.

(Ref: [Math156\\_classification\\_MNIST.ipynb](#) and [Math156\\_20NewsGroups.ipynb](#))

- (i) Reproduce Figure 16 for classifying digits ‘4’ and ‘7’ using multinomial naive Bayes. Compare the performance with logistic regression and probit regression.

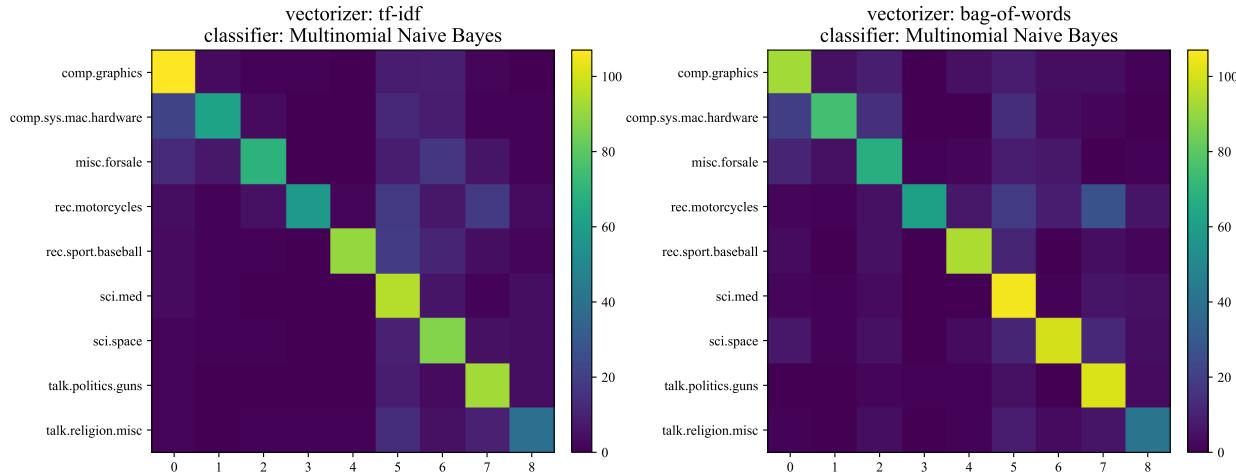


FIGURE 29. Confusion matrices for the 10-class document classification for the 20Newsgroups dataset with multinomial naive Bayes. Vectorizers of bag-of-words and tf-idf are both used.

- (ii) Reproduce Figure 19 for classifying digits 0 through 4 using multinomial naive Bayes. Compare the performance with multiclass logistic regression.

**Exercise 2.8.2.** In this exercise, we investigate the ‘robustness’ of multinomial naive Bayes (MNB) classifier on the 20Newsgroups dataset. (Ref: [Math156\\_20NewsGroups.ipynb](#))

- Modify the data preprocessing function “sample\_multiclass\_20NEWS” in the Jupyter notebook so that it has additional parameters ‘noise\_ratio’ and ‘noise\_intensity’: If ‘noise\_ratio’ is  $\epsilon \in [0, 1]$ , and ‘noise\_intensity’ is  $\delta > 0$ , then then it will add independent Uniform([0, 1]) variables times  $\delta$  to  $\epsilon$  density of randomly chosen coordinates.
- Choose two categories in the 20Newsgroups dataset and compare the performance of binary classification using MNB with varying noise parameters  $(\epsilon, \delta) \in \{0.1, 0.5, 0.9\} \times \{0.001, 0.01, 0.1\}$ . (Use tf-idf vectorizer.) Do you think if MNB is robust against noise on the 20Newsgroups dataset? If so, is it more robust against noise ratio ( $\epsilon$ ) or noise intensity ( $\delta$ )? Otherwise, can you say why?
- Perform similar experiments as in (ii) for multiclass logistic regression.

## 9. Gaussian Naive Bayes

Consider the case of classification problem where the features are continuous (e.g., height, weight, pixel values, tf-idf frequency) and may also assume negative values (e.g., bitcoin price change, EEG signal). In fact, in Exercise 2.7.5, we did discuss a naive Bayes model for continuous and nonnegative features, which was a straightforward generalization of the multinomial naive Bayes for nonnegative integer valued features. In this section, we introduce another widely used naive Bayes model called the *Gaussian naive Bayes* (GNB) model. The gist of it is the following: Assume each feature coordinate is a Gaussian RV with its own mean and variance, independently conditional on the class.

- Consider a  $\kappa$ -class classification problem for class labels  $\mathcal{Y} = \{1, \dots, \kappa\}$  and the feature vector  $\phi(\mathbf{x}) \in \mathbb{R}^M$  of each data point  $\mathbf{x}$  assume general *real* values in each coordinate (including negative values).
- We assume the following probabilistic model for  $(\text{data}, \text{class}) = (\mathbf{X}, Y)$ :

$$(\text{prior}) \quad \mathbb{P}(Y = i) =: q_i \quad \text{for } i \in \{1, \dots, \kappa\} \quad (169)$$

$$(\text{class-conditional PDF}) \quad p(\phi(\mathbf{X}) = [c_1, \dots, c_M] | Y = i) \propto \prod_{j=1}^M f_{\mu_{ij}, \sigma_{ij}}(c_j) \quad \text{for } i \in \{1, \dots, \kappa\}, \quad (170)$$

$$(\text{constraints}) \quad \begin{cases} 0 \leq q_1, \dots, q_\kappa \leq 1, & q_1 + \dots + q_\kappa = 1 \\ \sigma_{i,j} \geq 0 & \text{for all } i \in \{1, \dots, \kappa\}, j \in \{1, \dots, M\}, \end{cases} \quad (171)$$

where  $f_{\mu,\sigma}(x)$  is the PDF of  $N(\mu, \sigma)$ :

$$f_{\mu,\sigma}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right). \quad (172)$$

We denote  $\mathbf{W} = (q_i, \mu_{ij}, \sigma_{ij})$  for the vector of parameters of our model above.

Suppose training examples are given:  $(\phi(\mathbf{x}_s), y_s)$ ,  $s = 1, \dots, N$ , and write  $\phi(\mathbf{x}_s) = [\phi_1(\mathbf{x}_s), \dots, \phi_N(\mathbf{x}_s)]$  for each  $s$ . We appeal to the general maximum likelihood formulation of naive Bayes models in (131) and obtain the following:

$$\hat{\mathbf{W}} = \arg \max_{\mathbf{W}} \left[ \sum_{s=1}^N \sum_{i=1}^\kappa p \mathbf{1}(y_s = i) \log q_i + \sum_{s=1}^N \sum_{j=1}^p \sum_{i=1}^\kappa \mathbf{1}(y_s = i) \log \frac{1}{\sqrt{2\pi\sigma_{ij}^2}} \exp\left(-\frac{(\phi_j(\mathbf{x}_s) - \mu_{ij})^2}{2\sigma_{ij}^2}\right) \right] \quad (173)$$

$$= \arg \max_{\mathbf{W}} \left[ \ell_{GNB}(\mathbf{W}) := p \sum_{i=1}^\kappa \left( \sum_{s=1}^N \mathbf{1}(y_s = i) \right) \log q_i - \frac{1}{2} \sum_{j=1}^p \sum_{i=1}^\kappa \left( \sum_{s=1}^N \mathbf{1}(y_s = i) \right) \log(2\pi\sigma_{ij}^2) - \sum_{s=1}^N \sum_{j=1}^p \sum_{i=1}^\kappa \mathbf{1}(y_s = i) \frac{(\phi_j(\mathbf{x}_s) - \mu_{ij})^2}{2\sigma_{ij}^2} \right]. \quad (174)$$

The optimization problem (173) can be solved in closed form using standard arguments. The details are given in (2.9.1).

**Exercise 2.9.1.** Let  $\ell_{GNB}$  denote the Gaussian naive Bayes loss function defined in (173). Denote the solution of (173) as  $\hat{\mathbf{W}} = (\hat{\mathbf{q}}, \hat{\mu}_{ij}, \hat{\sigma}_{ij})$ .

(i) Argue that the optimal prior PMF  $\hat{\mathbf{q}} := [\hat{q}_1, \dots, \hat{q}_\kappa]$  solves the following optimization problem

$$\hat{\mathbf{q}} = \arg \max_{\substack{\mathbf{q} = [q_1, \dots, q_\kappa] \\ \text{PMF on } \{1, \dots, \kappa\}}} \sum_{i=1}^\kappa \left( \sum_{s=1}^N \mathbf{1}(y_s = i) \right) \log q_i. \quad (175)$$

Conclude that  $\hat{q}_i = \frac{1}{N} \sum_{s=1}^N \mathbf{1}(y_s = i)$  (Hint: Use Exercise 2.7.2).

(ii) Show that

$$\frac{\partial \ell_{GNB}(\mathbf{W})}{\partial \mu_{ij}} = \sum_{s=1}^N \mathbf{1}(y_s = i) \frac{(\phi_j(\mathbf{x}_s) - \mu_{ij})}{\sigma_{ij}^2}. \quad (176)$$

Deduce that  $\hat{\mu}_{ij}$  equals the following ‘sample mean of the  $j$ th feature in class  $i$ ’:

$$\hat{\mu}_{ij} = \frac{\sum_{s=1}^N \mathbf{1}(y_s = i) \phi_j(\mathbf{x}_s)}{\sum_{s=1}^N \mathbf{1}(y_s = i)}. \quad (177)$$

(Compare this with  $\hat{q}_{ij}$  in Prop. 2.7.1.)

(iii) Show that

$$\frac{\partial \ell_{GNB}(\mathbf{W})}{\partial \sigma_{ij}} = -\sigma_{ij}^{-1} \left( \sum_{s=1}^N \mathbf{1}(y_s = i) \right) + \sigma_{ij}^{-3} \sum_{s=1}^N \mathbf{1}(y_s = i) (\phi_j(\mathbf{x}_s) - \mu_{ij})^2. \quad (178)$$

Deduce that  $\hat{\sigma}_{ij}$  equals the following ‘variance of the class- $i$  empirical distribution’:

$$\hat{\sigma}_{ij}^2 = \frac{\sum_{s=1}^N \mathbf{1}(y_s = i) (\phi_j(\mathbf{x}_s) - \hat{\mu}_{ij})^2}{\sum_{s=1}^N \mathbf{1}(y_s = i)}, \quad (179)$$

where  $\hat{\mu}_{ij}$  is given in (ii).

## CHAPTER 3

# Neural Networks

In this chapter, we discuss a simple form of multilayer feedforward neural network. The model is simple enough to discuss without too much technical details but also is sophisticated enough to convey essential ideas in modern neural network models.

### 1. Model formulation

Very simply put, one can think of a feedword neural network as a composition of multiple general linear models (e.g., logistic regression, probit regression). To motivate this construction, recall that in generalized linear models, we model the predictive probabilities as

$$\hat{\mathbf{y}} := \begin{bmatrix} p(1 | \boldsymbol{\phi}(\mathbf{x})) \\ \vdots \\ p(\kappa | \boldsymbol{\phi}(\mathbf{x})) \end{bmatrix} = \sigma(\underbrace{\boldsymbol{\phi}(\mathbf{x})^T \mathbf{W}}_{\text{predictive prob.}}), \quad (180)$$

where  $\sigma$  is an activation function (e.g., signomid, logit, ReLu, tanh. See Figure 31.),  $\boldsymbol{\phi} \in \mathbb{R}^M$  is a  $M$ -dimensional feature vector of data  $\mathbf{x}$ ,  $\mathbf{W} \in \mathbb{R}^{M \times K}$  is a parameter matrix to be learned. So far, we have chosen a particular feature vector  $\boldsymbol{\phi}(\mathbf{x})$  before we fit the model to the training data (e.g., pixel values for MNIST images, bag-of-words or tf-idf vectors for 20Newgroups texts). But how do we know what feature should we choose? Could we somehow also *learn* what feature to use during the training process?<sup>1</sup>

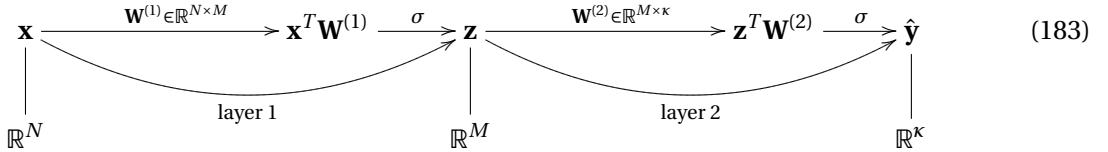
A simple idea is to *recycle the same generalized linear model in order to model the feature vector*. Namely, we consider

$$\mathbf{z} := \boldsymbol{\phi}(\mathbf{x}) = \sigma(\mathbf{x}^T \mathbf{W}'), \quad (181)$$

where we view  $\mathbf{x} \in \mathbb{R}^N$  and  $\mathbf{W}' \in \mathbb{R}^{N \times M}$  is another parameter matrix. Rewriting the parameter vectors as  $\mathbf{W}^{(1)} = \mathbf{W}' \in \mathbb{R}^N$  and  $\mathbf{W}^{(2)} = \mathbf{W} \in \mathbb{R}^M$ , the whole model reads as

$$\hat{\mathbf{y}} = \sigma(\sigma(\mathbf{x}^T \mathbf{W}^{(1)})^T \mathbf{W}^{(2)}). \quad (182)$$

As a diagram, it can also be described as (see also Figure 30)



This model is called a *two-layer neural network* (a.k.a. *multilayer perceptron*). Since the information only flows across layers in one direction, it is also called a *feedforward* neural network<sup>2</sup>. It is also called

<sup>1</sup>In the language of linear algebra, the feature vectors are coordinates on a chosen basis. We have worked with a fixed basis, but here we are attempting to also learn the basis from the training data. We will revisit this point later in dimensionality reduction (e.g., principal component analysis, nonnegative matrix factorization)

<sup>2</sup>Recurrent Neural Networks, which is a current state-of-the-art model for natural language processing, use feedback loop.

a *fully connected* neural network, since each coordinate in a given layer can get information from all coordinates in the previous layer<sup>3</sup>.

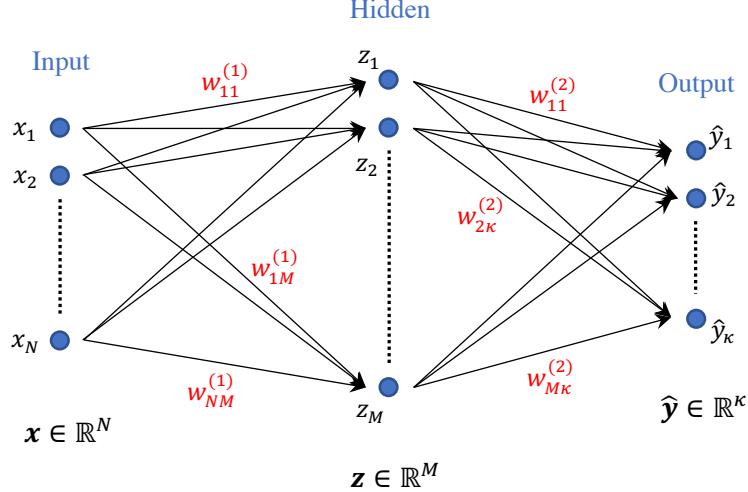


FIGURE 30. Diagrammatic representation of a 2-layer neural network.

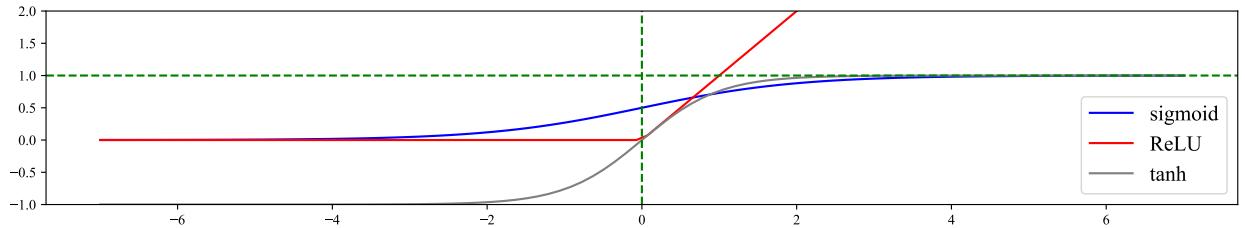


FIGURE 31. Common activation functions used in neural networks. Sigmoid:  $x \mapsto 1/(1 + e^{-x})$ , ReLU:  $x \mapsto \max(0, x)$ , and tanh:  $x \mapsto (e^{2x} - 1)/(e^{2x} + 1)$ .

## 2. Training Feedforward Neural Networks

In this section, we discuss how we train feedforward neural networks.

**2.1. Optimization framework.** As always, the key step is to formulate a suitable optimization problem for training the model. Suppose we are given training examples  $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)$  in  $\mathbb{R}^n \times \mathbb{R}^\kappa$ . Denote the predicted output for each input  $\mathbf{x}_s$  using parameter  $\mathbf{w}$  by  $\hat{\mathbf{y}}(\mathbf{x}_s; \mathbf{w}) = [\hat{y}_1(\mathbf{x}_s; \mathbf{w}), \dots, \hat{y}_\kappa(\mathbf{x}_s; \mathbf{w})]$ . We think of them as the output computed by the two-layer neural network (183), where  $\mathbf{w}$  denotes the collection of all parameter matrices  $[\mathbf{W}^{(1)}, \mathbf{W}^{(2)}]$ <sup>4</sup>. We seek to find an optimal parameter  $\hat{\mathbf{w}}$  that minimizes the overall training loss:

$$\hat{\mathbf{w}} := \arg \min_{\mathbf{w}} \left[ \ell_{NN}(\mathbf{w}) := \sum_{s=1}^n \ell(\mathbf{y}_s, \hat{\mathbf{y}}(\mathbf{x}_s; \mathbf{w})) \right], \quad (184)$$

where  $\ell$  is a loss function that measures the discrepancy between the true output  $\mathbf{y}_i$  and the predicted output  $\hat{\mathbf{y}}(\mathbf{x}_i; \mathbf{w})$ , which may depend on whether we are solving a regression or a classification problem<sup>5</sup>.

<sup>3</sup>Convolutional Neural Network, which is a current state-of-the-art for image classification, uses locally connected convolutional layers, mimicking the human visual cortex.

<sup>4</sup>However, our discussion in this section is general and not necessarily limited to the 2-layer neural networks.

<sup>5</sup>For classification problems, we may represent each class  $i$  by the indicator vector (one-hot encoding)  $[0, \dots, 0, 1, 0, \dots, 0] \in \mathbb{R}^\kappa$ , where the 1 is at the  $i$ th coordinate.

The standard choices are

$$\text{Regression} \quad (\text{square loss}): \quad \ell(\mathbf{y}_s, \hat{\mathbf{y}}(\mathbf{x}_s; \mathbf{w})) = \frac{1}{2} \|\mathbf{y}_s - \hat{\mathbf{y}}(\mathbf{x}_s; \mathbf{w})\|_F^2, \quad (185)$$

$$\text{Classification} \quad (\text{cross-entropy loss}): \quad \ell(\mathbf{y}_s, \hat{\mathbf{y}}(\mathbf{x}_s; \mathbf{w})) = - \sum_{i=1}^{\kappa} \mathbf{1}(\mathbf{y}_s = \text{class } i) \log \hat{y}_i(\mathbf{x}_s; \mathbf{w}), \quad (186)$$

$$(\text{softmax-cross-entropy loss}): \quad \ell(\mathbf{y}_s, \hat{\mathbf{y}}(\mathbf{x}_s; \mathbf{w})) = - \sum_{i=1}^{\kappa} \mathbf{1}(\mathbf{y}_s = \text{class } i) \log \text{softmax}_i(\hat{\mathbf{y}}(\mathbf{x}_s; \mathbf{w})). \quad (187)$$

The square loss for the regression case is natural and we have seen it for polynomial regression models in Section 2. On the other hand, for classification problems, recall that we usually model the predictive probabilities given features so that  $\hat{\mathbf{y}}(\mathbf{x}_s; \mathbf{w}) \in \mathbb{R}^\kappa$  can be viewed as a predictive PMF of the class of  $\mathbf{x}_s$ . One can think of the ‘cross-entropy’ as a similarity measure between two probability distributions. Namely, if we have two PMFs  $\mathbf{p} = [p_1, \dots, p_\kappa]$  and  $\mathbf{q} = [q_1, \dots, q_\kappa]$ , then the *cross-entropy* of  $\mathbf{q}$  relative to  $\mathbf{p}$  is defined by

$$H(\mathbf{p}, \mathbf{q}) = - \sum_{i=1}^{\kappa} p_i \log q_i. \quad (188)$$

Note that the following class-indicator vector

$$[\mathbf{1}(\mathbf{y}_s = \text{class } 1), \dots, \mathbf{1}(\mathbf{y}_s = \text{class } \kappa)] \quad (189)$$

is a PMF on the  $\kappa$  classes. Hence the cross-entropy loss in (186) is in fact the cross-entropy between the class-indicator vector above with the predictive PMF  $\hat{\mathbf{y}}(\mathbf{x}_s; \mathbf{w}) \in \mathbb{R}^\kappa$ . Lastly, the softmax-cross-entropy loss in (187) is frequently used in classification problems where the target output  $\mathbf{y}_s$  is a PMF but the predicted output  $\hat{\mathbf{y}}(\mathbf{x}_s; \mathbf{w})$  is not and is yet to be converted into a PMF. In this case, one can use the following softmax transform

$$\mathbf{x} = [x_1, \dots, x_\kappa]^T \rightarrow \text{softmax}(\mathbf{x}) = \left[ \frac{\exp(x_1)}{\sum_{i=1}^{\kappa} \exp(x_i)}, \dots, \frac{\exp(x_n)}{\sum_{i=1}^{\kappa} \exp(x_i)} \right]^T \in \mathbb{R}^\kappa \quad (190)$$

first and then compute the cross-entropy loss, which results in the combined loss of the form in (187).

One of the recurring theme in all of regression and classification models we studied so far is to take the *loss function as the negative log likelihood*. Namely,

$$\ell_{\text{total}}(\mathbf{w}) := -\log L(y_1, \dots, y_n; \mathbf{w}). \quad (191)$$

In fact, the square loss and cross-entropy loss above can be derived from this maximum likelihood formulation. For instance, consider the case of  $\kappa$ -class classification. Recall that  $\hat{\mathbf{y}}(\mathbf{x}_s; \mathbf{w})$  is the predictive PMF so that its  $i$ th coordinate,  $\hat{y}_i(\mathbf{x}_s; \mathbf{w})$ , is the probability that the class of  $\mathbf{x}_s$  is  $i$  according to our model. Hence the joint likelihood is

$$L(y_1, \dots, y_N; \mathbf{w}) = \mathbb{P}(Y_1 = y_1, \dots, Y_N = y_N; \mathbf{w}) = \prod_{s=1}^n \prod_{i=1}^{\kappa} (\hat{y}_i(\mathbf{x}_s; \mathbf{w}))^{\mathbf{1}(\mathbf{y}_s = \text{class } i)} \quad (192)$$

Hence we have

$$-\log L(y_1, \dots, y_N; \mathbf{w}) = - \sum_{s=1}^n \sum_{i=1}^{\kappa} \mathbf{1}(\mathbf{y}_s = \text{class } i) \log(\hat{y}_i(\mathbf{x}_s; \mathbf{w})), \quad (193)$$

which is exactly the total loss function in (184) using the cross-entropy loss in (186)<sup>6</sup>.

**Exercise 3.2.1.** Derive the total loss function in (184) using the square loss in (185) by using the maximum likelihood framework in (191). You may assume that the true output  $\mathbf{y} \in \mathbb{R}^\kappa$  is generated by a  $k$ -dimensional multivariate Normal distribution  $N(\hat{\mathbf{y}}(\mathbf{x}; \mathbf{w}), \sigma^2 I)$ , where  $I$  is the  $\kappa \times \kappa$  identity matrix. (See Example 3.2.2.)

---

<sup>6</sup>Compare this computation with the one for the multiclass logistic regression in Section 3.3.

**Example 3.2.2** (Multivariate Gaussian distribution). The a  $p$ -dimensional *multivariate Gaussian distribution* is a function  $\mathbb{R}^p \rightarrow \mathbb{R}$  given by

$$p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{p/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right], \quad (194)$$

where  $\boldsymbol{\mu} \in \mathbb{R}^p$  and  $\boldsymbol{\Sigma} \in \mathbb{R}^{p \times p}$  is a symmetric positive semi-definite matrix<sup>7</sup> and  $|\boldsymbol{\Sigma}|$  denotes the determinant of  $\boldsymbol{\Sigma}$ . If a random vector  $\mathbf{X} \in \mathbb{R}^p$  has distribution given by (194), then we denote  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . Then it has the following properties

$$\boldsymbol{\mu} = \mathbb{E}[\mathbf{X}], \quad \boldsymbol{\Sigma} = \text{Cov}(\mathbf{X}) = \mathbb{E}[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T]. \quad (195)$$

For these reasons, we call  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  the *mean* and the *covariance matrix* of  $\mathbf{X}$ , respectively. ▲

**2.2. Complexity of the optimization problem.** The first thing to note for solving the optimization problem (184) for the two-layer neural networks is that it is a *non-convex* optimization problem. Namely, in both cases of square and cross-entropy loss, the total loss  $\ell_{\text{total}}(\mathbf{w})$  depends on the model output  $\hat{\mathbf{y}}(\mathbf{x}_s; \mathbf{w})$  in a convex manner but the model output depends on the weight matrices  $\mathbf{w} = [\mathbf{W}^{(1)}, \mathbf{W}^{(2)}]$  in a non-convex manner. On the contrary, recall that the cross-entropy loss for logistic regression, or a single-layer neural network, is a convex function on the parameters (see Exercise 2.3.2). So, unfortunately, it is going to be very difficult (practically impossible) to find the global minimum of the total loss for training a 2-layer neural network.

Another nice way to see the complexity of the current optimization problem is using model symmetry. Observe that the model output of a 2-layer neural network (see Figure 30) is invariant under permuting the order of hidden units along with the corresponding in- and out-edge weights. For instance, if we swap  $z_1$  and  $z_N$ , then the weights  $w_{i1}^{(1)}$  from input  $x_i$  to  $z_1$  should be renamed to  $w_{iN}^{(1)}$  and also the weights  $w_{N1}^{(1)}$  from  $x_i$  to  $z_N$  should be renamed to  $w_{i1}^{(1)}$ , and so on; also the out-weights  $w_{1j}^{(2)}$  and  $w_{Nj}^{(2)}$  should be swapped. The model output  $\mathbf{y}$  should be unaware of this permutation since the values of the two linear transforms in the network is unchanged. Note that there are  $M!$  ways to permute the hidden units and their neural weights in this way. A consequence of this symmetry observation is that there are at least  $M!$  ‘equivalent local minima’ of the total loss function (184)<sup>8</sup>. For instance, if we use  $M = 10$  hidden units, then there are at least  $10! = 3,628,800$  global minima to look for.

Since the training of the 2-layer neural network has a complex ‘optimization landscape’, we change our goal from the ambitious one of finding a global minimum to the more modest one of finding a local minimum. It is a consistent experience in the literature that a local minimum found by some standard non-convex optimization algorithm (e.g., stochastic gradient descent, block coordinate descent) will perform as well as any other local minima, and this phenomenon belongs to an active area of machine learning research.

**2.3. Stochastic Gradient Descent and Backpropagation.** In order to find a local minimum of of (183), consider the familiar optimizatoin algorithm of gradient descent, which we introduced in discussing training of logistic regression in (82):

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta_t \nabla_{\mathbf{w}} \ell_{NN}(\mathbf{w}_t) \quad (196)$$

One of the main bottleneck in the development of neural network models in machine learning literature was the huge computational cost in computing the full gradient  $\ell_{NN}(\mathbf{w}_t)$  for using the above gradient descent algorithm. We will discuss two ways to overcome this issue: 1) Stochastic Gradient Descent (SGD) and 2) Backpropagation.

<sup>7</sup>A matrix  $\mathbf{A} \in \mathbb{R}^{p \times p}$  is *symmetric* if  $\mathbf{A} = \mathbf{A}^T$  and *positive semi-definite* if  $\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$  for all  $\mathbf{x} \in \mathbb{R}^p$ . Moreover,  $\mathbf{A}$  is *positive definite* if it is positive semi-definite and furthermore,  $\mathbf{x}^T \mathbf{A} \mathbf{x} > 0$  if  $\mathbf{x} \neq \mathbf{0}$ .

<sup>8</sup>For the hyperbolic tangent activation function  $\sigma(x) = \tanh(x) = (e^{2x} - 1)/(e^{2x} + 1)$ , the model output is also invariant under sign flips of the neural weights, so there are at least  $M!2^M$  equivalent local minimia.

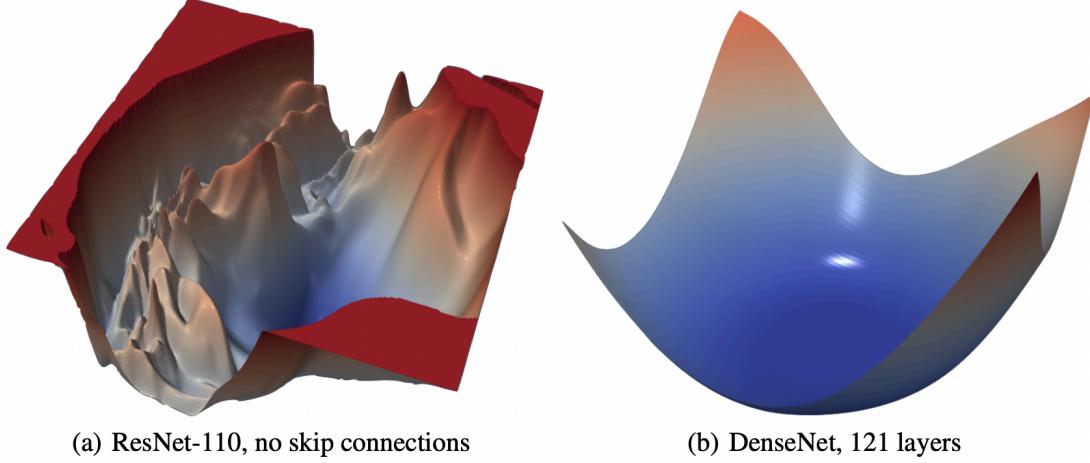


FIGURE 32. Optimization landscape for ResNet-110-noshort (left, nonconvex) and DenseNet for CIFAR-10 (right, convex). Figure excerpted from [LXT<sup>+</sup>17].

**2.3.1. Stochastic Gradient Descent.** The idea of SGD is not limited only to the context of neural networks, but is applicable for general nonconvex optimization problems. Simply speaking it is to use only a ‘partial gradient’ in each iteration of the gradient descent, which typically requires much less computational cost. In our case, note that we can decompose the gradient  $\nabla \ell_{NN}(\mathbf{w}_t)$  as the sum of gradients per sample:

$$\nabla \ell_{NN}(\mathbf{w}_t) = \sum_{s=1}^n \nabla_{\mathbf{w}} \ell(\mathbf{y}_s, \hat{\mathbf{y}}(\mathbf{x}_s; \mathbf{w}_t)), \quad (197)$$

by using the definition of total loss function  $\ell_{NN}$  and the linearity of differentiation. Hence, instead of computing the full gradient  $\nabla \ell_{\mathbf{w}} \ell(\mathbf{w}_t)$ , why don’t we choose a small set  $I \subseteq \{1, \dots, n\}$  of training data samples and compute the corresponding partial gradient

$$\ell_{NN; I} \ell(\mathbf{w}_t) := \sum_{i \in I} \nabla_{\mathbf{w}} \ell(\mathbf{y}_s, \hat{\mathbf{y}}(\mathbf{x}_s; \mathbf{w}_t)). \quad (198)$$

The stochastic component of SGD lies in choosing the index set  $I$  randomly. Hence the corresponding SGD algorithm takes the following form:

---

**Algorithm 1** Stochastic Gradient Descent for Neural Networks

---

- 1: **Input:**  $\mathbf{w}_0$  (initial parameter);  $T$  (number of iterations);  $(\eta_t)_{t \geq 1}$ , (non-increasing weights in  $(0, 1]$ );  $n_0 \in \{1, \dots, n\}$  (size of subsample)
  - 2:   **for**  $t = 1, \dots, T$  **do**:
  - 3:     Sample  $I \subseteq \{1, \dots, n\}$  with  $|I| = n_0$  over all subsets of size  $n_0$ ;
  - 4:      $\mathbf{w}_t \leftarrow \mathbf{w}_{t-1} - \eta_t \nabla_{\mathbf{w}} \ell_{NN; I} \ell(\mathbf{w}_{t-1})$
  - 5:   **end for**
  - 6: **output:**  $\mathbf{w}_T$
- 

An extreme case of the SGD algorithm in Algorithm 1 is when the size  $n_0$  of subsample is 1. In that case, in each iteration, we only use a single randomly chosen training example to update the current parameter  $\mathbf{w}_t$ . In general, the chosen set of training examples  $\{\mathbf{x}_i ; i \in I\}$  is called a *minibatch*, and SGD amounts to process only a randomly selected minibatch in each iteration of gradient descent, instead of the full *batch*  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ . Processing minibatches instead of the full batch is what reduces the computational cost of SGD. Moreover, the randomized choice of minibatches is what helps the SGD not to overfit to some particular training examples and to escape from undesirable local minima.

It is a classical result in the theory of nonconvex optimization that the parameters  $(\mathbf{w}_t)_{t \geq 0}$  produced by SGD converges to a stationary point of the loss function (183) if a ‘diminishing step size’ is used:  $\sum_{t=1}^{\infty} \eta_t = \infty$  and  $\sum_{t=1}^{\infty} \eta_t^2 < \infty$  when the loss function is nice enough (see, e.g., [SSY18]). As in the case of logistic regression, we will use the schedule  $\eta_t = \gamma(\log t)/\sqrt{t}$ , where  $\gamma > 0$  is a constant.

**2.3.2. Backpropagation.** *Error backpropagation* is an effective method to *evaluate* the gradient of a feedforward neural network at the current parameters, first developed by Rumelhart, Hinton, and Williams [RHW86].

In order to simplify the notation, we only consider the loss at  $s$ th training sample  $(\mathbf{x}_s, \mathbf{y}_s)$ :

$$\ell_s(\mathbf{w}) := \ell(\mathbf{y}_s, \hat{\mathbf{y}}(\mathbf{x}_s; \mathbf{w})). \quad (199)$$

In order to execute the SGD algorithm for training neural networks, we need to compute the following derivatives: For each  $s \in \{1, \dots, n\}$ ,

$$\frac{\partial \ell_s}{\partial w_{ij}^{(1)}}, \quad \frac{\partial \ell_s}{\partial w_{jq}^{(2)}}. \quad (200)$$

Namely, we want to differentiate the loss function with respect to the edge weights. However, it would be extremely helpful for us to compute the derivative of the loss function with respect to *node weights*  $z_j$  and  $\hat{y}_q$ . Recall the following relations:

$$z_j = \sigma(\mathbf{x}^T \mathbf{W}^{(1)}[:, j]) = \sigma\left(\sum_{i=1}^N x_i w_{ij}^{(1)}\right), \quad \hat{y}_q = \sigma(\mathbf{z}^T \mathbf{W}^{(2)}[:, q]) = \sigma\left(\sum_{j=1}^M z_j w_{jq}^{(2)}\right). \quad (201)$$

Denoting  $\mathbf{y}_s = [(\mathbf{y}_s)_1, \dots, (\mathbf{y}_s)_K]$ , we first differentiate the loss function with respect to  $\hat{y}_q$ :

$$(\text{square loss}): \quad \frac{\partial \ell_s}{\partial \hat{y}_q} = \frac{1}{2} \frac{\partial}{\partial \hat{y}_q} \|\mathbf{y}_s - \hat{\mathbf{y}}(\mathbf{x}_s; \mathbf{w})\|_F^2 = (\hat{y}_q - (\mathbf{y}_s)_q), \quad (202)$$

$$(\text{cross-entropy loss}): \quad \frac{\partial \ell_s}{\partial \hat{y}_q} = -\frac{\partial}{\partial \hat{y}_q} \sum_{i=1}^K \mathbf{1}(\mathbf{y}_s = \text{class } i) \log \hat{y}_i(\mathbf{x}_s; \mathbf{w}) = -\frac{\mathbf{1}(\mathbf{y}_s = \text{class } q)}{\hat{y}_q}. \quad (203)$$

$$(\text{softmax-cross-entropy loss}): \quad \frac{\partial \ell_s}{\partial \hat{y}_q} = -\frac{\partial}{\partial \hat{y}_q} \sum_{i=1}^K \mathbf{1}(\mathbf{y}_s = \text{class } i) \log \text{softmax}_i(\hat{\mathbf{y}}(\mathbf{x}_s; \mathbf{w})) \quad (204)$$

$$= \text{softmax}_q(\hat{\mathbf{y}}(\mathbf{x}_s; \mathbf{w})) - (\mathbf{y}_s)_q. \quad (205)$$

For (205), see Exercise 3.2.3.

Next, we compute the derivatives of the loss function with respect to the hidden state  $z_j$ . By using chain rule:

$$\underbrace{\frac{\partial \ell_s}{\partial z_j}}_{\text{'error' at node } z_j} = \sum_{q=1}^K \underbrace{\frac{\partial \ell_s}{\partial \hat{y}_q} \frac{\partial \hat{y}_q}{\partial z_j}}_{\text{'error' at node } \hat{y}_q} = \sum_{q=1}^K \underbrace{\frac{\partial \ell_s}{\partial \hat{y}_q}}_{\text{'error' at node } \hat{y}_q} \sigma'\left(\underbrace{\mathbf{z}^T \mathbf{W}^{(2)}[:, q]}_{\text{activation for } \hat{y}_q}\right) w_{jq}^{(2)}, \quad (206)$$

where  $\sigma'$  denotes the derivative of the activation function  $\sigma$ . The equation in (206) is an instance of error backpropagation (See Figure 35).

In order to explain the mechanics of error backpropagation, we introduce the notion of ‘errors at nodes’ of a neural network to be the loss function with respect to the state of the node<sup>9</sup>. For instance,  $\frac{\partial \ell_s}{\partial \hat{y}_q}$  is the error at node  $\hat{y}_q$  and  $\frac{\partial \ell_s}{\partial z_j}$  is the error at node  $z_j$ , and so on<sup>10</sup>. With this notion at hand, in (202) and (203), we computed the errors at the output nodes  $\hat{y}_j$  directly. Also, we can view (206) as a *recursive relation* that computes the error at node  $z_j$  in terms of the errors at nodes  $\hat{y}_q$  in the output layer. Below is a verbal description of the error backpropagation for 2-layer neural networks:

<sup>9</sup>We are making a slight abuse of notation of identifying a node in a neural network with its state.

<sup>10</sup>Why such a derivative is called an ‘error’? One can make sense of it at least for the square loss at the output layer. See (202).

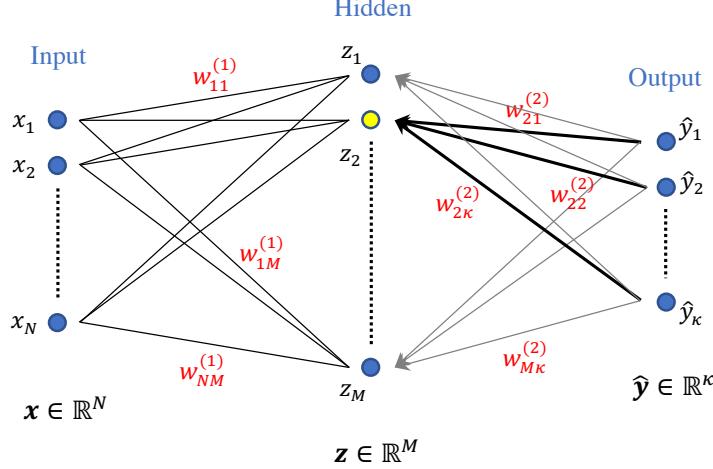


FIGURE 33. Error backpropagation in a 2-layer neural network. The ‘errors’ computed at the output layer are backpropagated to the second hidden node through the edges, where the corresponding edge weight is multiplied. The total sum of such weighted backpropagated errors at node \$z\_2\$ is then the ‘error’ at node \$z\_2\$.

### Computation of gradients for training a 2-layer FFNN

1. Forward propagate the input all the way to the output layer and compute the activation and state of every node.
2. Compute the errors at the output nodes using (202)-(205).
3. Backpropagate the errors in the output layer to the hidden layer and compute the errors at hidden nodes, using (206).

In order to see that such recursion can be used to compute the errors in deeper layers, it is instrumental to compute the errors in the input layer, even though we do not use them for training the 2-layer neural network. One can easily verify

$$\underbrace{\frac{\partial \ell_s}{\partial x_i}}_{\text{‘error at node } x_i} = \sum_{j=1}^M \frac{\partial \ell_s}{\partial z_j} \frac{\partial z_j}{\partial x_i} = \sum_{j=1}^M \underbrace{\frac{\partial \ell_s}{\partial z_j}}_{\text{‘error’ at } z_j} \sigma' \left( \underbrace{\mathbf{x}^T \mathbf{W}^{(1)}[:, j]}_{\text{activation for } z_j} \right) w_{ij}^{(1)}. \quad (207)$$

Hence, if we had a deep neural network with 2 or more layers, then we could use the same backpropagation recursion layer by layer to compute the errors at each node in all layers.

Now that we can compute the errors at each node, it remains to see how we use them to compute the desired derivatives of the loss function with respect to the neural weights in (200). This can be easily done by the relation (201) and chain rule:

$$\frac{\partial \ell_s}{\partial w_{jq}^{(2)}} = \frac{\partial \ell_s}{\partial \hat{y}_q} \frac{\partial \hat{y}_q}{\partial w_{jq}^{(2)}} = \frac{\partial \ell_s}{\partial \hat{y}_q} \sigma'(\mathbf{z}^T \mathbf{W}^{(2)}[:, q]) z_j, \quad (208)$$

$$\frac{\partial \ell_s}{\partial w_{ij}^{(1)}} = \frac{\partial \ell_s}{\partial z_j} \frac{\partial z_j}{\partial w_{ij}^{(1)}} = \frac{\partial \ell_s}{\partial z_j} \sigma'(\mathbf{x}^T \mathbf{W}^{(1)}[:, j]) x_i. \quad (209)$$

Lastly, with all the derivatives as above, we can now execute the SGD algorithm (Algorithm 1) by aggregating the above per-sample derivatives to get the minibatch gradient \$\nabla\_{\mathbf{w}} \ell\_{NN;I}(\mathbf{w}\_{t-1})\$ by using (198).

**Exercise 3.2.3** (Derivative of softmax-cross-entropy). The softmax function  $\sigma : \mathbb{R}^n \rightarrow \mathbb{R}^n$ , is defined by

$$\mathbf{x} = [x_1, \dots, x_n]^T \mapsto \sigma(\mathbf{x}) = [\sigma_1(\mathbf{x}), \dots, \sigma_n(\mathbf{x})] = \left[ \frac{\exp(x_1)}{\sum_{i=1}^n \exp(x_i)}, \dots, \frac{\exp(x_n)}{\sum_{i=1}^n \exp(x_i)} \right]^T. \quad (210)$$

Define the following softmax-cross-entropy loss

$$\ell(\mathbf{y}, \hat{\mathbf{y}}) := - \sum_{i=1}^K \mathbf{1}(\mathbf{y} = \text{class } i) \log \sigma_i(\hat{\mathbf{y}}) \quad (211)$$

for each  $\mathbf{y}, \hat{\mathbf{y}} \in \mathbb{R}^n$ .

(i) Show that

$$\frac{\partial}{\partial x_j} \sigma_i(\mathbf{x}) = \sigma_i(\mathbf{x}) (\mathbf{1}(i = j) - \sigma_j(\mathbf{x})). \quad (212)$$

(ii) Show that

$$\frac{\partial}{\partial \hat{\mathbf{y}}} \ell(\mathbf{y}, \hat{\mathbf{y}}) = \sigma(\hat{\mathbf{y}}) - \mathbf{y}. \quad (213)$$

**Exercise 3.2.4.** Formulate an algorithm that computes the gradient of a general  $L$ -layer neural network. You may need to extend the error backpropagation algorithm for 2-layer NN we discussed above to the general  $L$ -layer case and show its correctness.

### 3. Examples of 2-layer FFNN

**Example 3.3.1.** Consider a two-layer neural network with the following choices:

*Loss function:*  $\ell_s(\mathbf{y}_s, \hat{\mathbf{y}}(\mathbf{x}_s; \mathbf{w})) = (1/2) \|\mathbf{y}_s - \hat{\mathbf{y}}(\mathbf{x}_s; \mathbf{w})\|_F^2$  (Square loss)

*Activation functions:*  $\tanh(x) = (e^x - e^{-x}) / (e^x + e^{-x})$  for the hidden layer and sigmoid  $\sigma(x)$  for the output layer.

Then the total loss function is given by

$$\nabla \ell_{NN}(\mathbf{w}) = \frac{1}{2} \sum_{s=1}^n \|\mathbf{y}_s - \hat{\mathbf{y}}(\mathbf{x}_s; \mathbf{w})\|_F^2. \quad (214)$$

After forward-propagating an input  $\mathbf{x}_s \in \mathbb{R}^N$ , the hidden states  $\mathbf{z}_s \in \mathbb{R}^M$  and the output  $\hat{\mathbf{y}}_s = \hat{\mathbf{y}}(\mathbf{x}_s; \mathbf{w}) \in \mathbb{R}^K$  are given by

$$\mathbf{z}_s = \tanh(\mathbf{x}_s^T \mathbf{W}^{(1)}), \quad \hat{\mathbf{y}}_s = \sigma(\mathbf{z}_s^T \mathbf{W}^{(2)}). \quad (215)$$

For backpropagation, the errors at the output nodes are given by

$$\frac{\partial \ell_s}{\partial (\hat{\mathbf{y}}_s)_q} = ((\hat{\mathbf{y}}_s)_q - (\mathbf{y}_s)_q), \quad (216)$$

Hence, from (206), the errors at the hidden nodes are given by

$$\frac{\partial \ell_s}{\partial (\mathbf{z}_s)_j} = \sum_{q=1}^K \frac{\partial \ell_s}{\partial (\hat{\mathbf{y}}_s)_q} ((\hat{\mathbf{y}}_s)_q (1 - (\hat{\mathbf{y}}_s)_q)) w_{jq}^{(2)}, \quad (217)$$

where we have used that  $\sigma'(x) = \sigma(x)(1 - \sigma(x))$ .

Now from (208), we get

$$\frac{\partial \ell_s}{\partial w_{jq}^{(2)}} = \frac{\partial \ell_s}{\partial (\hat{\mathbf{y}}_s)_q} ((\hat{\mathbf{y}}_s)_q (1 - (\hat{\mathbf{y}}_s)_q)) (\mathbf{z}_s)_j. \quad (218)$$

On the other hand, note that  $\tanh'(x) = 1 - \tanh^2(x)$  so

$$\tanh'(\mathbf{x}_s^T \mathbf{W}^{(1)}[:, j]) = 1 - [(\mathbf{z}_s)_j]^2. \quad (219)$$

Then (209) gives

$$\frac{\partial \ell_s}{\partial w_{ij}^{(1)}} = \frac{\partial \ell_s}{\partial (\mathbf{z}_s)_j} (1 - [(\mathbf{z}_s)_j]^2) (\mathbf{x}_s)_i. \quad (220)$$

Lastly, the derivatives of the total loss  $\ell_{NN}$  with respect to the neural weights are given by summing up the per-sample derivatives above over  $s$  in the selected minibatch  $I \subseteq \{1, \dots, n\}$ .  $\blacktriangle$

**Exercise 3.3.2** (2-layer FFNN on MNIST). Below we show the train and test accuracy of classifying MNIST dataset (digits 0 through 4) using the 2-layer FFNN discussed in Example 3.3.1. Here the accuracy for multiclass classification is computed by the total number of correctly classified examples divided by the total number of examples. (Ref: [Jupyter Notebook](#))

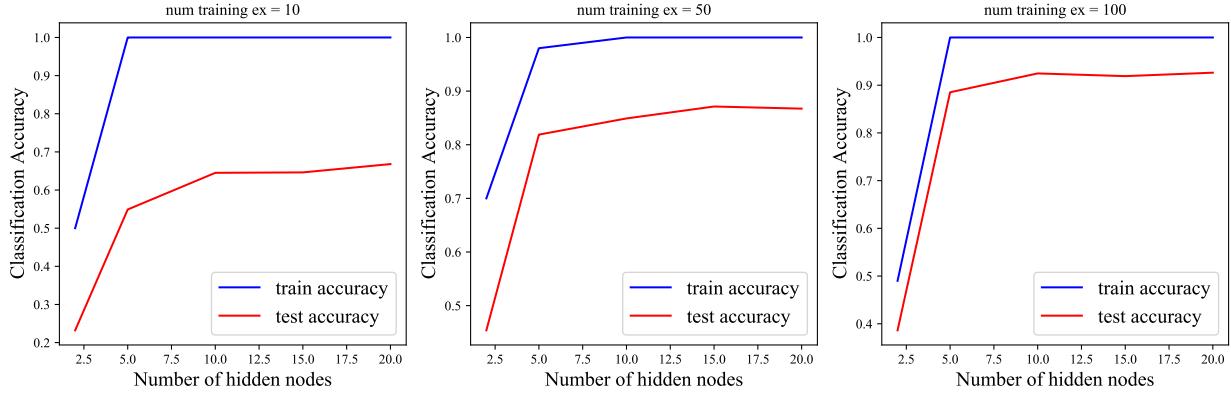


FIGURE 34. Classification accuracy of 2-layer FFNN for MNIST data with digits from 0 up to 4. The number of hidden nodes are varied within  $\{2, 5, 10, 15, 20\}$  and the number of training examples are varied within  $\{10, 50, 100\}$ .

- (i) What can you conclude from Figure 34? Optimal number of hidden nodes for each size of training set? Overfitting in terms of the number of hidden nodes and the size of training set?
- (ii) Make similar plots for (1) digits 5 through 9 and (2) all digits while keeping other parameters the same. (You may increase the number of iterations for training for classifying all digits.) Do your observations in (ii) still hold in these experiments?

**Exercise 3.3.3.** Consider the following  $L_2$ -regularized 2-layer neural network training problem:

$$\hat{\mathbf{w}} := \underset{\mathbf{w} = [\mathbf{W}^{(1)}, \mathbf{W}^{(2)}]}{\operatorname{argmin}} \left[ \ell_{NN}(\mathbf{w}; \lambda) := \left( \sum_{s=1}^n \ell(\mathbf{y}_s, \hat{\mathbf{y}}(\mathbf{x}_s; \mathbf{w})) \right) + \lambda (\|\mathbf{W}^{(1)}\|_F^2 + \|\mathbf{W}^{(2)}\|_F^2) \right], \quad (221)$$

where  $\lambda \geq 0$  is a  $L_2$ -regularization constant. This will penalize the weight matrices  $\mathbf{W}^{(1)}, \mathbf{W}^{(2)}$  to have large  $L_2$ -norm.

- (i) How does this change the training algorithm? Derive the gradient of the above loss function and formulate the modified training algorithm.
- (ii) Modify the Multilayer Perceptron implementation MLP in [Jupyter Notebook](#) so that it has the option to use  $L_2$ -regularization. Reproduce Figure 34 with varying levels of regularization. Report a reasonable value of the regularization constant for the experiment in Figure 34.

**Exercise 3.3.4.** Implement the two-layer neural network discussed in Example 3.3.1 on your own, where the activation function in the output layer should allow the identity function in order for applying it on regression problems. Using your code, reproduce a figure similar to Figure 35 below. For each  $n = 6, 10$  data points in the sinusoidal data set (see (15)), what do you think is the best value of  $M$ ?

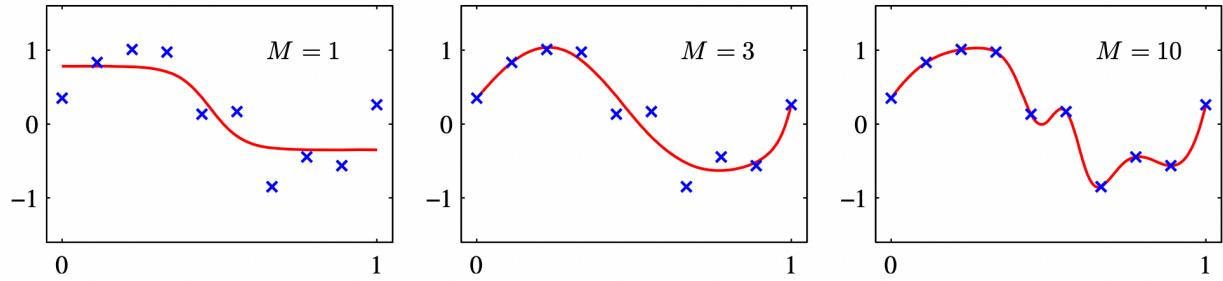


FIGURE 35. Examples of two-layer networks trained on 10 data points drawn from the sinusoidal data set (see Sec. 2.2). The graphs show the result of fitting networks having  $M = 1, 3$  and 10 hidden units, respectively, by minimizing a sum-of-squares error function using a SGD algorithm. Figure excerpted from [Bis06].

#### 4. Convolutional Neural Networks

**4.1. Introduction.** In this section, we discuss one of the most successful and architecture of feed-forward deep neural networks for image classification, namely, the *convolutional neural networks* (CNN, or ConvNet). A few highlighting features of CNN are as follow:

1. CNN is a ‘localized’ version of the fully connected deep feedforward neural network (inspired by how animal visual cortex works).
2. CNN learns patterns in images that are ‘shift-invariant’, meaning that the location of a certain geometric pattern in the image matters less<sup>11</sup>.
3. CNN not only to recognize patterns from images, but also to abstract such patterns. (e.g., first layer learns edges or colors, next layer learns larger curves and shapes, and so on.)

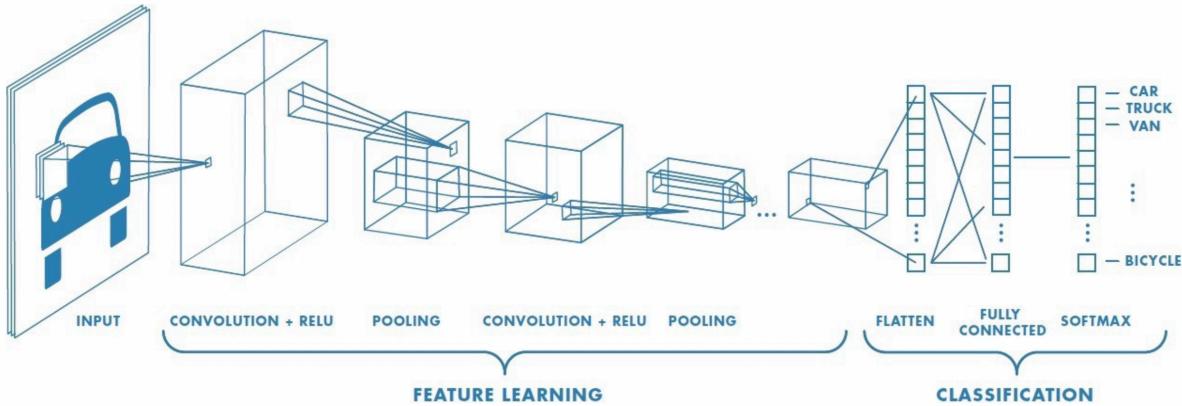


FIGURE 36. Architecture of a Convolutional Neural Network for image classification. Figure by S. Saha [link](#)

We first discuss a motivating example. Consider the now-familiar MNIST handwritten digits image dataset. We have seen that regardless of what algorithm we use for classifying the digits, we usually get very high accuracy (e.g., see Figure 16, 17, 22 for binary classification with LR and PR, Figure 19 for multiclass classification with MLR). Since the images in the MNIST dataset are so well aligned within  $28 \times 28$  frame, it only means that we can classify the images to the correct classes against variation within each class under such a good alignment. So what if the images are not perfectly aligned within the same frame, as in Figure 37? There, we ‘padded’ additional zeros around the original image to perturb the

<sup>11</sup>Think about whether the usual FFNN or multiclass logistic regression have such abilities

alignment. In other words, we embedded each  $28 \times 28$  image inside  $48 \times 48$  zero matrices at uniformly randomly chosen locations. Do you think the classification algorithms we have developed so far would still work well after this perturbation?<sup>12</sup>

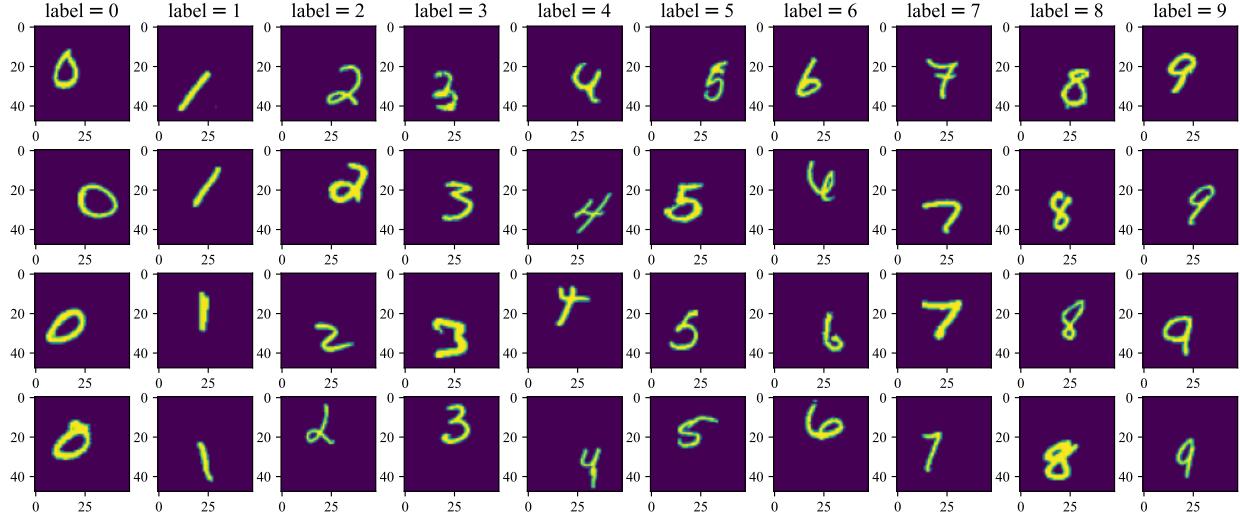


FIGURE 37. Examples of MNIST dataset of handwritten digits by Yann LeCun, Corinna Cortes, and Christopher J.C. Burges with random padding. Each  $28 \times 28$  image is superimposed within a  $48 \times 48$  zero matrix at uniformly chosen random location.

In Figure 38, we show the training and testing accuracies of using 2-layer neural network with 10 hidden nodes for classifying MNIST (five digits from 0 through 4) with random padding. We vary the ‘thickness’ of the padding, namely, when the thickness is  $r$ , it means that each  $28 \times 28$  image is randomly embedded inside a zero matrix of shape  $(28 + r) \times (28 + r)$ . As we can see in Figure 38, the training accuracies are close to 1 but the testing accuracies drop quickly as we increase the padding thickness. This means our model overfits to some non-shift-invariant features. Can we overcome this shortcoming by increasing the number of layers in our model? Not quite, as we see in Figure 39 for 3-layer FFNN.

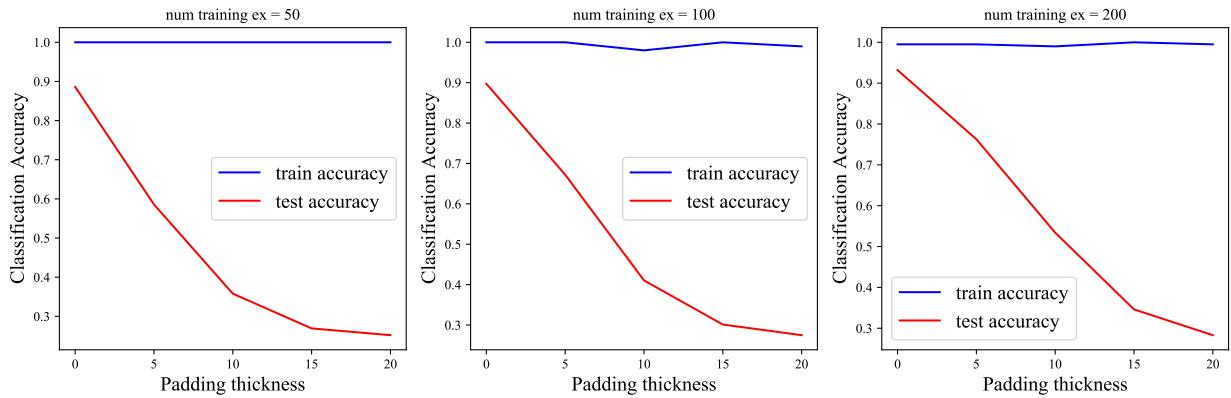


FIGURE 38. Classification accuracy of 2-layer FFNN for MNIST data with digits from 0 up to 4 with random padding with varying thickness. The number of hidden nodes is 10. When the thickness is  $r$ , it means that each  $28 \times 28$  image is randomly embedded inside a zero matrix of shape  $(28 + r) \times (28 + r)$ .

<sup>12</sup>Certainly our brain doesn't care the location of the digits within a larger image.

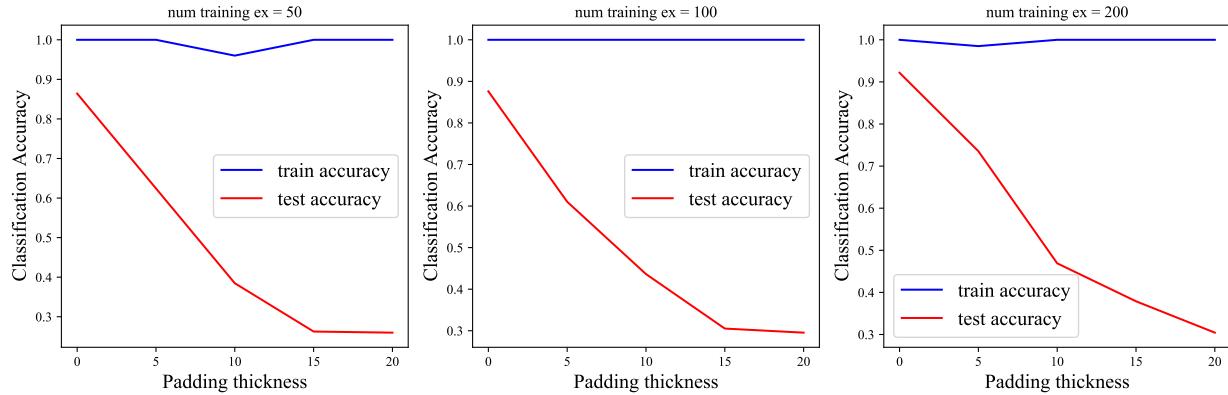


FIGURE 39. Classification accuracy of 3-layer FFNN for MNIST data with digits from 0 up to 4 with random padding with varying thickness. The number of hidden nodes are 10 for each hidden layer. When the thickness is  $r$ , it means that each  $28 \times 28$  image is randomly embedded inside a zero matrix of shape  $(28 + r) \times (28 + r)$ .

**Exercise 3.4.1.** Reproduce figures similar to Figure 38 using (1) multiclass logistic regression; and (2) naive Bayes classifier. Conclude that whether these algorithms are able to learn shift-invariant features. (Ref: [Jupyter Notebook – ‘Classifying non-aligned MNIST images’](#))

**4.2. Model formulation.** In order to formulate CNN, we need to introduce (1) convolutional layer and (2) pooling layer and how they work. Each ‘layer’ in a CNN is the following composition: convolutional with set of filters + ReLU + pool. Every time we complete such procedure, we learn more abstract and shift-invariant features from the previous layer. (See Figure 40 for an illustration<sup>13</sup>)

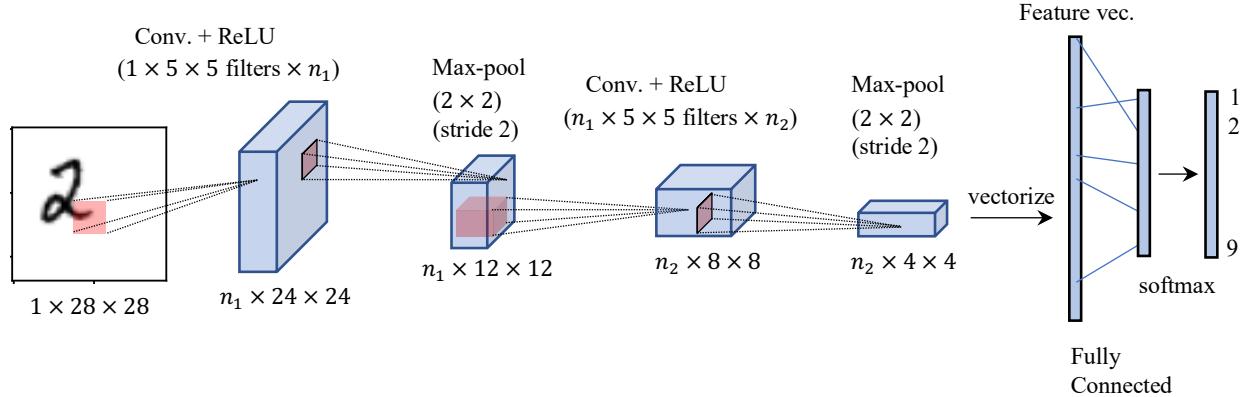


FIGURE 40. An example of CNN architecture for MNIST image classification. Convolution + ReLU activation followed by a max-pool layers are applied twice with to extract shift-invariant feature vector. Pass this into a multiclass logistic regression (fully connected layer + softmax) for classification.

**4.2.1. Convolutional layer.** First we explain how a convolutional layer works in CNN. Recall that in a fully connected neural network, the previous layer’s activation (node state), say  $\mathbf{z}$ , is dot producted with each columns of the following weight matrix to produce the activation of each nodes in the following layer. In CNN, it is useful to note that

$$\text{Weight matrix in a CNN layer} = \text{Set of tensors (filters/kernels)} F_1, \dots, F_m \in \mathbb{R}^{c \times r \times r}, \quad (222)$$

<sup>13</sup>See also ‘A Comprehensive Guide to Convolutional Neural Networks’ by S. Saha

where the ‘tensor depth’  $c$  should match the depth of the previous ‘image’. For example, color images have shape  $3 \times x \times y$  and monochrome images have shape  $1 \times x \times y$ . Each filter is supposed to capture a particular pattern *everywhere* in the image. We explain this in steps as follow:

1. Say we have a *filter*  $F$  of shape  $1 \times 3 \times 3$ . For an input  $X$  of shape  $1 \times 3 \times 3$ , the convolution  $X * F$  of  $X$  with filter  $F$  is the sum of entrywise product of  $X$  with  $F$ .

$$\underbrace{\begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{bmatrix}}_X * \underbrace{\begin{bmatrix} f_{11} & f_{12} & f_{13} \\ f_{21} & f_{22} & f_{23} \\ f_{31} & f_{32} & f_{33} \end{bmatrix}}_F = \text{sum} \left( \begin{bmatrix} x_{11}f_{11} & x_{12}f_{12} & x_{13}f_{13} \\ x_{21}f_{21} & x_{22}f_{22} & x_{23}f_{23} \\ x_{31}f_{31} & x_{32}f_{32} & x_{33}f_{33} \end{bmatrix} \right) = \sum_{i,j} x_{ij}f_{ij}. \quad (223)$$

For general depth  $c$ , the convolution of a filter  $F$  of shape  $c \times r \times r$  on an input of shape  $1 \times r \times r$  is a scalar given similarly as follow:

$$X * F = \sum_{i,j,k} X[i, j, k] F[i, j, k], \quad (224)$$

where the sum is over all entries  $(i, j, k)$ .

2. If the input  $X$  is  $1 \times 4 \times 4$ , then we can choose four  $1 \times 3 \times 3$  *patches*<sup>14</sup> inside  $X$  by *sliding*<sup>15</sup> a window (frame) of  $1 \times 3 \times 3$ . We can convolve each patch with  $F$  to get a scalar. So output of  $X * F$  is a  $1 \times 2 \times 2$  tensor.

$$\underbrace{\begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ x_{31} & x_{32} & x_{33} & x_{34} \\ x_{41} & x_{42} & x_{43} & x_{44} \end{bmatrix}}_X * \underbrace{\begin{bmatrix} f_{11} & f_{12} & f_{13} \\ f_{21} & f_{22} & f_{23} \\ f_{31} & f_{32} & f_{33} \end{bmatrix}}_F = \underbrace{\begin{bmatrix} X[0:3, 0:3] * F & X[0:3, 1:] * F \\ X[1:, 0:3] * F & X[1:, 1:] * F \end{bmatrix}}_{\text{In python slicing}}. \quad (225)$$

3. If we have a  $c \times r \times r$  filter  $F$  and  $c \times a \times b$  input  $X$ , then the result of convolution  $X * F$  with *stride* 1 is a  $1 \times (a - r + 1) \times (b - r + 1)$  tensor:

$$(X * F)[0, j, k] = X[:, i:i+r, j:j+r] * F. \quad (226)$$

Namely, the  $(i, j)$  entry of  $X * F$  is convolution of the  $(i, j)$  patch  $X[:, i:i+r, j:j+r] \in \mathbb{R}^{c \times r \times r}$  with  $F$ .

4. If we have a  $c \times r \times r$  filter  $F$  and  $c \times a \times b$  input  $X$ , then the result of convolution  $X * F$  with *stride*  $s$  is a  $1 \times (\frac{a-r}{s} + 1) \times (\frac{b-r}{s} + 1)$  tensor:

$$(X * F)[0, j, k] = X[:, is:is+r, js:js+r] * F. \quad (227)$$

Namely, the  $(i, j)$  entry of  $X * F$  is convolution of the  $(i, j)$  patch  $X[:, is:is+r, js:js+r] \in \mathbb{R}^{c \times r \times r}$  with  $F$ .

4. The result of the convolutional layer with input  $X \in \mathbb{R}^{c \times a \times b}$  and filters  $F_1, \dots, F_m \in \mathbb{R}^{c \times r \times r}$  with stride  $s$  is

$$X * (\text{filters } F_1, \dots, F_m) = [X * F_1, \dots, X * F_m] \in \mathbb{R}^{m \times (\frac{a-r}{s} + 1) \times (\frac{b-r}{s} + 1)}, \quad (228)$$

which we think of as a ‘depth’  $m$  image of size  $(\frac{a-r}{s} + 1) \times (\frac{b-r}{s} + 1)$ .

See Figure 40 for an illustration.

<sup>14</sup>Square sub-matrices of consecutive rows and columns

<sup>15</sup>We may want to skip every some rows and columns in extracting patches from  $X$  to avoid too much overlap between patches and/or reduce the output of convolution more rapidly. The number of rows/columns we slide over is called *stride*. For *stride*=1, we sample the maximum number of patches.

**4.2.2. Pooling layer.** In CNN, a *pooling* is used to reduce the size of the previous layer's output by either keeping the most dominant cell value in each patch (called 'max-pool') or averaging every patch into a single scalar (called 'avg-pool'). We introduce these concepts with the following examples.

1.  $2 \times 2$  pooling on  $1 \times 2 \times 2$  input  $X$ :

$$\text{pool}\left(X = \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{bmatrix}\right) = \begin{cases} \max(X) = \max\{x_{ij} \mid 1 \leq i, j \leq 2\} & \text{for max-pool} \\ \text{avg}(X) = (1/4) \sum_{1 \leq i, j \leq 2} x_{ij} & \text{for avg-pool} \end{cases}. \quad (229)$$

2.  $2 \times 2$  pooling on  $1 \times 4 \times 4$  input  $X$  with stride 2:

$$\text{pool}\left(X = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ x_{31} & x_{32} & x_{33} & x_{34} \\ x_{41} & x_{42} & x_{43} & x_{44} \end{bmatrix}\right) = \begin{bmatrix} \text{pool}\left(\begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{bmatrix}\right) & \text{pool}\left(\begin{bmatrix} x_{13} & x_{14} \\ x_{23} & x_{24} \end{bmatrix}\right) \\ \text{pool}\left(\begin{bmatrix} x_{31} & x_{32} \\ x_{41} & x_{42} \end{bmatrix}\right) & \text{pool}\left(\begin{bmatrix} x_{33} & x_{34} \\ x_{43} & x_{44} \end{bmatrix}\right) \end{bmatrix} \in \mathbb{R}^{1 \times 2 \times 2}. \quad (230)$$

3.  $2 \times 2$  pooling on  $1 \times 4 \times 4$  input  $X$  with stride 1:

$$\text{pool}\left(X = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ x_{31} & x_{32} & x_{33} & x_{34} \\ x_{41} & x_{42} & x_{43} & x_{44} \end{bmatrix}\right) = [\text{pool}(X[i:i+2, j:j+2] ; 0 \leq i, j \leq 2)] \in \mathbb{R}^{1 \times 3 \times 3}. \quad (231)$$

4.  $2 \times 2$  pooling on  $c \times a \times b$  input  $X$  with stride 1:

$$\text{pool}(X) = [\text{pool}(X[d, i:i+2, j:j+2] ; 0 \leq d < c, 0 \leq i \leq a-2, 0 \leq j \leq b-2)] \quad (232)$$

$$\in \mathbb{R}^{c \times (a-1) \times (b-1)}. \quad (233)$$

5.  $r \times r$  pooling on  $c \times a \times b$  input  $X$  with stride  $s$ :

$$\text{pool}(X) = [\text{pool}(X[d, is:is+r, j:j+r] ; 0 \leq d < c, 0 \leq i \leq \frac{a-r}{s}, 0 \leq j \leq \frac{b-r}{s})] \quad (234)$$

$$\in \mathbb{R}^{c \times (\frac{a-r}{s}+1) \times (\frac{b-r}{s}+1)}. \quad (235)$$

**Exercise 3.4.2.** Consider the CNN architecture shown in Figure 40. Suppose we use the following general choices of convolutional and pooling layers:

Conv1:  $n_1$  filters of shape  $1 \times r \times r$  with stride  $s_f$

Conv2:  $n_2$  filters of shape  $n_1 \times r \times r$  with stride  $s_f$

maxpool (both):  $d \times d$  maxpool with stride  $s_p$

Also assume that input image has dimension  $1 \times a \times b$ . Then compute the dimension of the feature vector that goes into the first dense layer as a function of  $n_1, n_2, r, d, s_f, s_p, a, b$ .

**4.3. Training.** Training a CNN is similar to that for the feedforward neural networks. Namely, we use stochastic gradient descent in order to find the optimal parameters (that is, all filters, weight matrices, and biases). In order to compute the gradients of the parameters, we first backpropagate a chosen training image, compute the error at the output layer, and then backpropagate it all the way down to the input layer. This will give the gradients of the loss with respect to node states, which we can use to compute the gradients of filters and weight matrices by chain rule.

In order to execute backpropagation for CNN, we need to perform the following four kinds of back-propagations:

1. Backprop through a dense layer — Same as in FFNN
2. Backprop through a maxpool layer — Assign errors to the 'winning pixels'
3. Backprop through a ReLU layer — Set the errors zero at coordinates where the input activation is negative
4. Backprop through a convolutional layer — Details omitted. (Ref: [Backprop. in CNN](#))

### Computation of gradients for training a CNN

1. Forward propagate the input all the way to the output layer and compute the activation and state of every node.
2. Compute the errors at the output nodes using (202)-(205).
3. Backpropagate the errors in the output layer to the hidden layer and compute the errors at hidden nodes, using (206).

**4.4. Example: MNIST + padding.** Next, we give some examples of using CNN for MNIST image classification. We are going to use MNIST with random padding as shown in Figure 37. The architecture of CNN we use here is as follow:

$$\text{Input} \rightarrow \underbrace{\text{Conv1}}_{\text{filter } 1 \in \mathbb{R}^{n_1 \times 1 \times 5 \times 5}} + \text{ReLU} + \text{maxpool} \rightarrow \underbrace{\text{Conv2}}_{\text{filter } 1 \in \mathbb{R}^{n_2 \times n_1 \times 5 \times 5}} + \text{ReLU} + \text{maxpool} \quad (236)$$

$$\rightarrow \underbrace{\text{Dense1} + \text{ReLU}}_{W1} \rightarrow \underbrace{\text{Dense2} + \text{Softmax}}_{W2} \quad (237)$$

$$\rightarrow \text{Output} \quad (238)$$

The parameters we need to optimize are  $n_1$  filters of shape  $1 \times 5$  in the first convolutional layer,  $n_2$  filters of shape  $n_1 \times 5 \times 5$  in the second convolutional layer, and the two weight matrices for the dense layers. For all pooling layer, we use  $2 \times 2$  maxpool with stride 2. All convolutions use stride 1.

In Figure 41, we show the training and testing accuracies of using the above CNN in (238) for classifying MNIST (five digits from 0 through 4) with random padding. Again we vary the thickness of the padding, namely, when the thickness is  $r$ , it means that each  $28 \times 28$  image is randomly embedded inside a zero matrix of shape  $(28+r) \times (28+r)$ . The overall trend is similar as for FFNN in Figure 38: The training accuracies are close to 1 but the testing accuracies drop quickly as we increase the padding thickness. However, the test accuracy is about 10% higher with 100 training examples and 20% higher with 200 training examples than those for FFNN. This means that the above CNN is able to learn shift-invariant features to some extent, at least much better than FFNN.

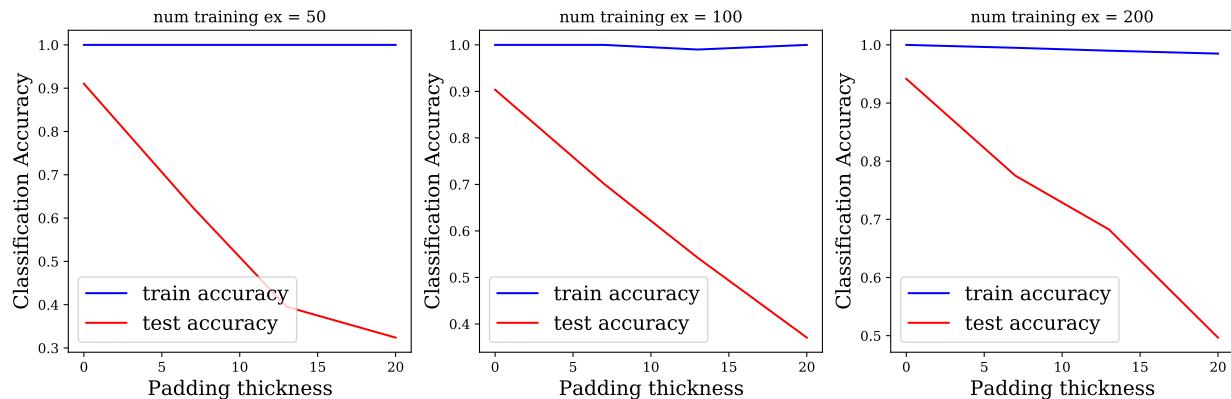


FIGURE 41. Classification accuracy of 3-layer FFNN for MNIST data with digits from 0 up to 4 with random padding with varying thickness. The number of hidden nodes are 10 for each hidden layer. When the thickness is  $r$ , it means that each  $28 \times 28$  image is randomly embedded inside a zero matrix of shape  $(28+r) \times (28+r)$ .

**Exercise 3.4.3.** Reproduce Figures 39 and 41 with larger number of training examples: 300, 500, 1000. Compare the performance of FFNN and CNN using your plots. (Ref: Jupyter notebook on [FFNN](#) and on [CNN](#))

**Exercise 3.4.4.** Produce a plot similar to Figure 34 for the CNN model in (238), where the horizontal axis correspond to the shared number  $n_1 = n_2 \in \{2, 5, 10, 20\}$  of filters in each convolutional layer. Use

padding thickness 10 and keep all the other settings the same as in Figure 41. What do you think the optimal number of filters to use is? (Ref: Jupyter notebook on [CNN](#))

In Figure 42, we show 8 filters of shape  $1 \times 5 \times 5$  in the first convolutional layer that are learned for a binary MNIST classification (digits 0 and 1). Can you imagine that kind of pattern that each filter is suppose to pick up from the image?

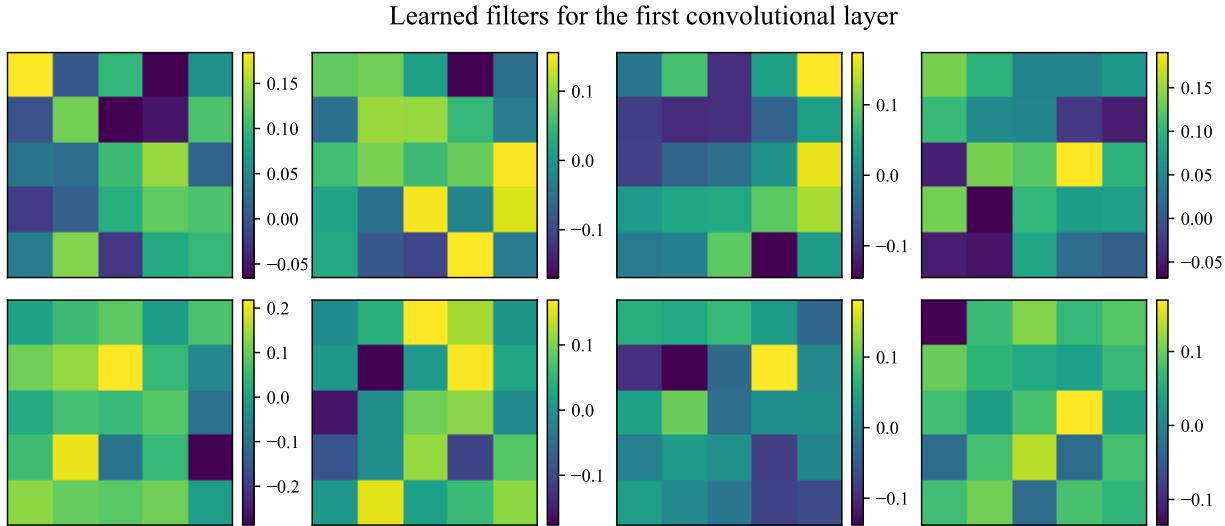


FIGURE 42. 8 filters learned for a binary MNIST classification (digits 0 and 1) using the CNN model in (238).

**Exercise 3.4.5.** Train the CNN in (238) for MNIST classification with digits 0 and 1. Pick two images of each digit. Using the trained model, plot the image after the first and the second convolutional layers for each of the four chosen images. (Ref: Jupyter notebook on [CNN](#)) (Hint: Modify the ‘predict’ function in `src/CNN.py` so that it also returns outputs of each convolutional layer.)

## CHAPTER 4

# Matrix Factorization and Dictionary Learning

In this chapter, we introduce the problem of matrix factorization, which is a powerful unifying framework for a number of unsupervised machine learning problems such as dimension reduction or dictionary learning. Two main techniques we will be learning are *nonnegative matrix factorization* (NMF) and *principal component analysis* (PCA).

### 1. Introduction

**1.1. Dictionary Learning.** *Dictionary-learning* algorithms are machine-learning techniques that are able to learn interpretable latent structures of complex data sets and are applied regularly in the data analysis of text and images. Such algorithms usually consist of two steps. First, one samples a large number of structured subsets of a data set (e.g., square patches of an image or collections of a few sentences of a text); Second, one finds a set of basis elements such that taking a *nonnegative* linear combination of them can successfully approximate each of the sampled mesoscale patches. Such a set of basis elements is called a *dictionary*, and one can interpret each basis element as a latent structure of the data set. As an example, consider the image of the artwork CYCLE by M. C. Escher in Figure 43a. We first sample 10,000 square patches of  $21 \times 21$  pixels, and we then use a nonnegative matrix factorization (NMF) [LS99] algorithm to find a dictionary with  $r = 25$  square patches (see Figure 43a). Each element of the learned dictionary describes a latent shape in the image.

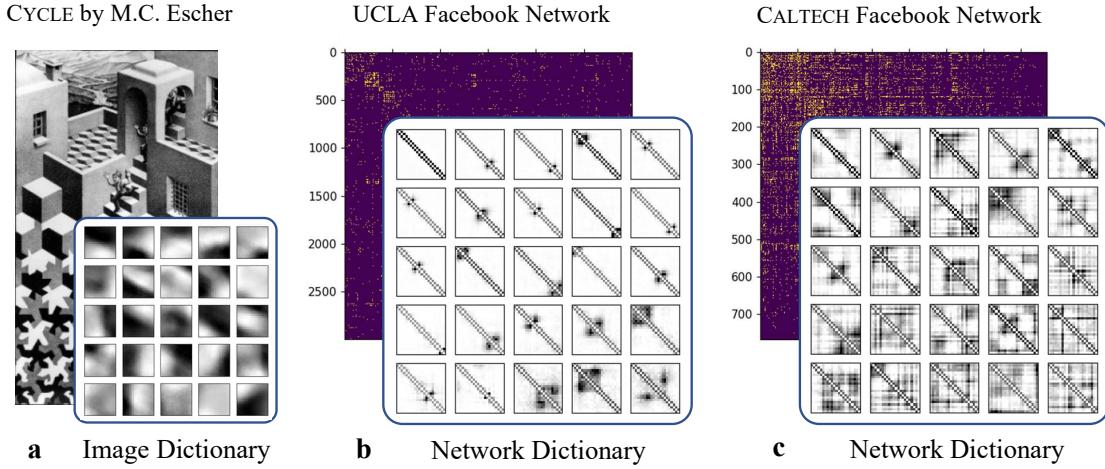


FIGURE 43. Illustration of matrix factorization. Each column of the data matrix is approximated by the linear combination of the columns of the dictionary matrix with coefficients given by the corresponding column of the code matrix. Figure excerpted from [LKVP21].

**1.2. Matrix Factorization.** *Matrix factorization* is one of the fundamental tools in dictionary learning problems. Given a large data matrix  $\mathbf{X}$ , can we find some small number of “dictionary vectors” so that we can represent each column of the data matrix as a linear combination of dictionary vectors? More precisely, given a data matrix  $\mathbf{X} \in \mathbb{R}^{d \times n}$  and a rank parameter  $r \in \mathbb{N}$ , we wish to factorize  $\mathbf{X}$  into the

product of  $\mathbf{W} \in \mathbb{R}^{d \times r}$  and  $\mathbf{H} \in \mathbb{R}^{r \times n}$  by solving the following optimization problem

$$(Matrix\ Factorization) \quad \inf_{\mathbf{W} \in \mathbb{R}^{d \times r}, \mathbf{H} \in \mathbb{R}^{r \times n}} \|\mathbf{X} - \mathbf{WH}\|_F^2. \quad (239)$$

Here  $\mathbf{W}$  is called the *dictionary* and  $\mathbf{H}$  is the *code* of data  $\mathbf{X}$  using dictionary  $\mathbf{W}$ . A solution of such matrix factorization problem is illustrated in Figure 44.

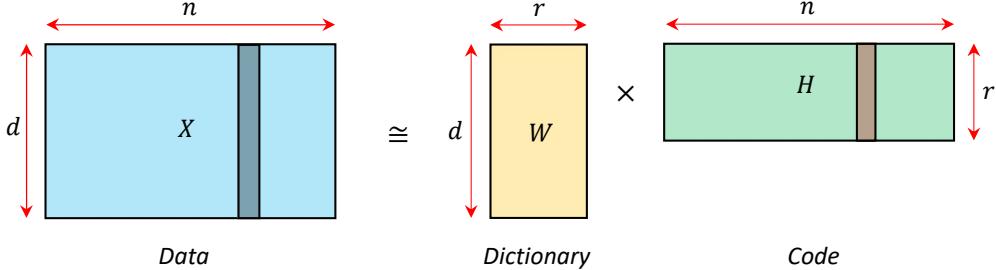


FIGURE 44. Illustration of matrix factorization. Each column of the data matrix is approximated by the linear combination of the columns of the dictionary matrix with coefficients given by the corresponding column of the code matrix. Figure excerpted from [LNB20].

Recall the method of least squared we studied for polynomial regression: (See also Exercise 1.2.1).

$$\hat{\mathbf{H}} = \underset{\mathbf{H} \in \mathbb{R}^{r \times n}}{\operatorname{argmin}} \underbrace{\|\mathbf{X} - \mathbf{W} \mathbf{H}\|_F^2}_{\text{given}} + \lambda \underbrace{\|\mathbf{H}\|_F^2}_{\text{TBD}}. \quad (240)$$

One can think of (240) as a matrix factorization problem where the first factor (dictionary) is known and we are only looking for the code  $\mathbf{H}$ . Because of this reason, we also call (240) as a *coding* problem, which is also called a *sparse coding* problem if the regularizer  $\lambda$  is positive<sup>1</sup>.

While the optimization problem 239 is an *unconstrained* problem in the sense that we are free to choose any dictionary  $\mathbf{W} \in \mathbb{R}^{d \times r}$  and code  $\mathbf{H} \in \mathbb{R}^{r \times n}$ , We will be more interested in its *constrained* versions, in particular using nonnegativity constraints. Consider the following modifications:

$$(Dictionary\ Learning) \quad \inf_{\mathbf{W} \in \mathbb{R}^{d \times r}, \mathbf{H} \in \mathbb{R}_{\geq 0}^{r \times n}} \|\mathbf{X} - \mathbf{WH}\|_F^2, \quad (241)$$

$$(NMF) \quad \inf_{\mathbf{W} \in \mathbb{R}_{\geq 0}^{d \times r}, \mathbf{H} \in \mathbb{R}_{\geq 0}^{r \times n}} \|\mathbf{X} - \mathbf{WH}\|_F^2. \quad (242)$$

In (241), we imposed an additional constraint that the code  $\mathbf{H} \in \mathbb{R}^{r \times n}$  should only take nonnegative entries. This allows us to interpret the columns of dictionary  $\mathbf{W}$  as ‘features’ that directly contribute to composing the columns of  $\mathbf{X}$  (samples)<sup>2</sup>. We can go even further and also impose the same nonnegativity constraint on the dictionary matrix  $\mathbf{W} \in \mathbb{R}^{r \times n}$  as in (242). This problem is called *nonnegative matrix factorization* (NMF), which is an extremely popular technique in modern machine learning for its high interpretability and mathematical simplicity.

**1.3. Applications in Image processing.** Before we discuss how to solve matrix factorization problems, we give some examples of matrix factorization in the context of image processing through image dictionary learning.

Suppose we are given an image, and we are interested in the following *unsupervised problem* of learning ‘essential features in all of its  $k \times k$  patches’. In other words, we would like to know what important features are there that if we cut out a large number of  $k \times k$  square patches from the original image.

<sup>1</sup>It is more effective to use an  $L_1$ -regularizer  $\lambda \|\mathbf{H}\|_1$  for the purpose of sparse coding. This is also closely related to (Lagrangian version of) the popular algorithm of LASSO in statistics.

<sup>2</sup>A general formulation of dictionary learning may also impose a convex constraint  $\mathcal{C} \subseteq \mathbb{R}^{d \times r}$  for dictionary  $\mathbf{W}$  along with the nonnegativity constraint on the code  $\mathbf{H}$ .

Of course, still this problem is not precisely defined. Even though we can state a precise optimization problem for it, let's instead describe the algorithm first and then see what exactly we are solving.

As we have briefly mentioned in Section (1.1), for dictionary learning we first do patch sampling and then feature extraction. For patch sampling, we get a large number  $N$  (e.g.,  $N = 10K$ ) of  $k \times k$  (e.g.,  $21 \times 21$ ) square patches from the original image. We will vectorize each such patch as a  $k^2$  dimensional column vector and form a data matrix  $\mathbf{X}$  of shape  $k^2 \times N$ . We then use an algorithm<sup>3</sup> for NMF to get a dictionary matrix  $\mathbf{W} \in \mathbb{R}^{k^2 \times r}$  and code matrix  $\mathbf{H} \in \mathbb{R}^{r \times N}$  with  $r = 25$ . Each of the 25 columns of  $\mathbf{W}$  are vectorized versions of the extracted features or ‘basis shapes’. Reshaping these  $k^2$  dimensional column vectors into  $k \times k$  square form, we obtain the image dictionary of size 25, as shown in Figure 45.

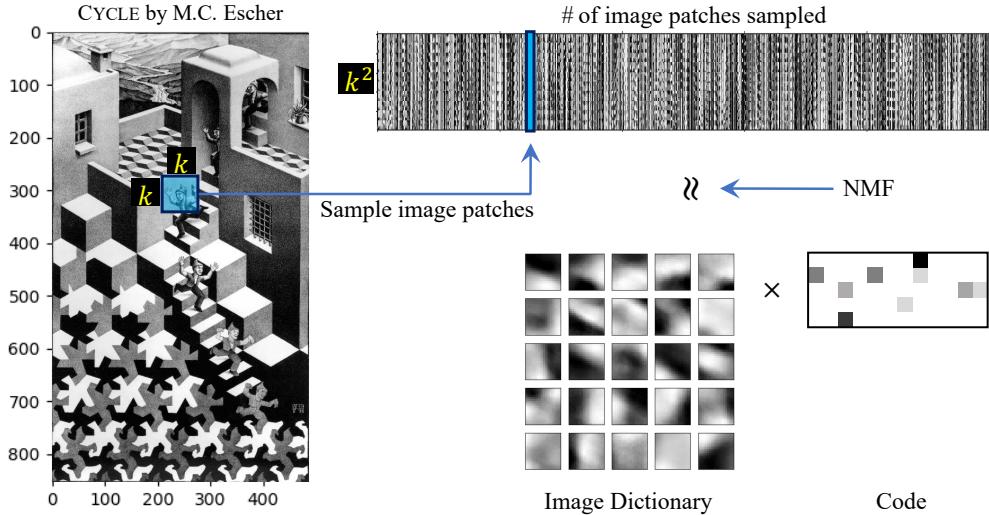


FIGURE 45. Examples of image dictionary learning by patch sampling and NMF.

After we learn a set of dictionary images from a given image, we can use them to reconstruct the original image. Moreover, if the given image is corrupted with some noise, then this reconstruction using learned image dictionary can be used for image denoising. Figure 46 gives an example of image dictionary learning and image denoising with noise separation. We omit details here.

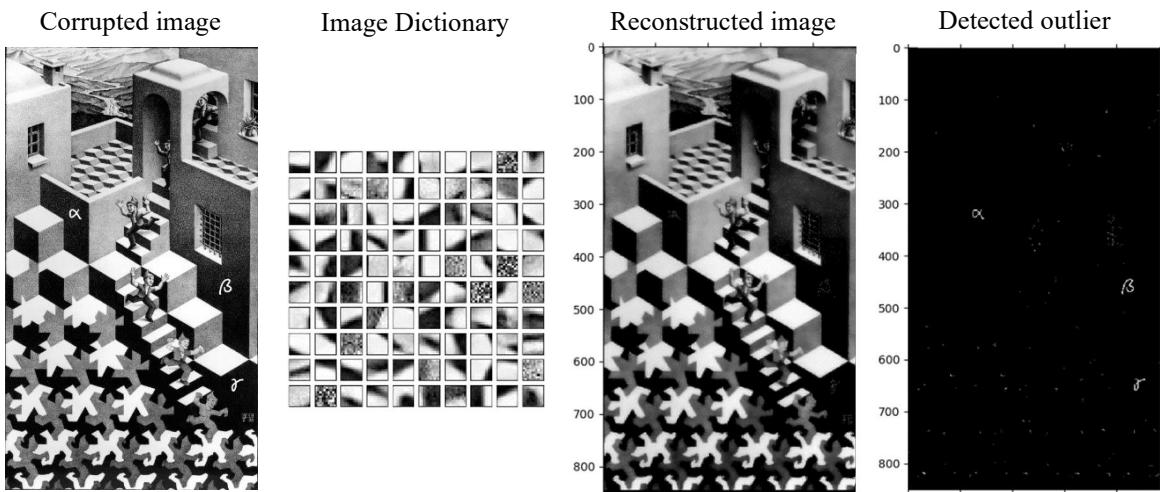


FIGURE 46. Examples of image dictionary learning by patch sampling and NMF. The learned image dictionary is used for denoising the original image corrupted with noises ‘ $\alpha$ ’, ‘ $\beta$ ’, and ‘ $\gamma$ ’.

<sup>3</sup>We will discuss algorithms for solving matrix factorization problems later

## 2. Algorithms for Matrix Factorization

In this section, we discuss algorithms for solving matrix factorization problems. First, fix a data matrix  $\mathbf{X} \in \mathbb{R}^{d \times n}$  and integer  $r \geq 1$ . Define the following loss function

$$\ell(\mathbf{W}, \mathbf{H}) := \|\mathbf{X} - \mathbf{WH}\|_F^2. \quad (243)$$

We would like to minimize the above loss function jointly in  $\mathbf{W} \in \mathbb{R}^{d \times r}$  and  $\mathbf{H} \in \mathbb{R}^{r \times n}$ , possibly with convex constraints (e.g., nonnegativity) depending on the situation.

**Exercise 4.2.1.** Show the following.

- (i) Show that the function  $\mathbf{H} \mapsto \|\mathbf{X} - \mathbf{WH}\|_F^2$  is convex.
- (ii) Show that the function  $\mathbf{W} \mapsto \|\mathbf{X} - \mathbf{WH}\|_F^2 = \|\mathbf{X}^T - \mathbf{H}^T \mathbf{W}^T\|_F^2$  is convex.
- (iii) Show that the function  $(\mathbf{W}, \mathbf{H}) \mapsto \|\mathbf{X} - \mathbf{WH}\|_F^2$  is non-convex.

According to Exercise 4.2.1, this is a non-convex optimization problem so it is difficult to find its global minima. However, if we fix one of the factors, then  $\ell$  is convex with respect to the other factor (see Exercise 4.2.1), which is in fact a least squares problem (see Exercise 1.2.1). Hence one can iteratively fix one factor and optimize the other, each time solving a least squares problem. This is called *Alternating Least Squares* (ALS) algorithm for matrix factorization, which we state in Algorithm 2 for NMF.

---

### Algorithm 2 Alternating Least Squares for Matrix Factorization

---

- 1: **Input:**  $\mathbf{X} \in \mathbb{R}_{\geq 0}^{d \times n}$  (data matrix);  $r \geq 1$  (rank parameter);  $[\mathbf{W}_0, \mathbf{H}_0] \in \mathbb{R}_{\geq 0}^{d \times r} \times \mathbb{R}_{\geq 0}^{r \times n}$  (initial estimate);  $N$  (number of iterations)
  - 2:   **for**  $n = 1, \dots, N$  **do**:
  - 3:      $\mathbf{H}_n \leftarrow \operatorname{argmin}_{\mathbf{H} \in \mathbb{R}_{\geq 0}^{r \times n}} \|\mathbf{X} - \mathbf{W}_{t-1} \mathbf{H}\|_F^2$
  - 4:      $\mathbf{W}_n \leftarrow \operatorname{argmin}_{\mathbf{W} \in \mathbb{R}_{\geq 0}^{d \times r}} \|\mathbf{X} - \mathbf{W} \mathbf{H}_n\|_F^2$
  - 5:   **end for**
  - 6: **output:**  $[\mathbf{W}_N, \mathbf{H}_N]$ .
- 

The ALS algorithm in Algorithm 2 is a special case of *block coordinate descent* (BCD) algorithms, which is to iteratively optimize the objective function with respect to a block coordinate (group of individual coordinates). The general algorithm is stated in Algorithm 3. For more backgrounds and its convergence properties, see [Lyu20].

---

### Algorithm 3 Block Coordinate Descent

---

- 1: **Requirement:**  $\boldsymbol{\theta} = [\theta_1, \dots, \theta_m] \mapsto f(\theta_1, \dots, \theta_m)$  is convex in each block-coordinate  $\theta_i$ ,  $i = 1, \dots, m$ .
  - 2: **Input:**  $\boldsymbol{\theta}_0 = (\theta_0^{(1)}, \dots, \theta_0^{(m)}) \in \Theta^{(1)} \times \dots \times \Theta^{(m)}$  (initial estimate);  $N$  (number of iterations)
  - 3:   **for**  $n = 1, \dots, N$  **do**:
  - 4:     Update estimate  $\boldsymbol{\theta}_n = [\theta_n^{(1)}, \dots, \theta_n^{(m)}]$  by
  - 5:       **For**  $i = 1, \dots, m$  **do**:
  - 6:          $\theta_n^{(i)} \in \operatorname{argmin}_{\theta \in \Theta^{(i)}} f(\theta_n^{(1)}, \dots, \theta_n^{(i-1)}, \theta, \theta_{n-1}^{(i+1)}, \dots, \theta_{n-1}^{(m)})$ ; (244)
  - 7:       **end for**
  - 8:   **end for**
  - 9: **output:**  $\boldsymbol{\theta}_N$
- 

**Exercise 4.2.2** (Nonnegative Coding with  $L_1$  regularizer). Suppose we have matrices  $\mathbf{X} \in \mathbb{R}^{d \times n}$  and  $\mathbf{W} \in \mathbb{R}^{d \times r}$ . We seek to find a matrix  $\hat{\mathbf{H}} \in \mathbb{R}_{\geq 0}^{r \times n}$  where

$$\hat{\mathbf{H}} = \operatorname{argmin}_{\mathbf{H} \in \mathbb{R}_{\geq 0}^{r \times n}} (\ell(\mathbf{H}) := \|\mathbf{X} - \mathbf{WH}\|_F^2 + \lambda_1 \|\mathbf{H}\|_1 + \lambda_2 \|\mathbf{H}\|_F^2). \quad (245)$$

Here  $\lambda \geq 0$  is called the  $L_1$ -regularization parameter. (This is an instance of *constrained* quadratic optimization problem.) (c.f. The use of  $L_1$  regularizer instead of  $L_2$  as in Exercise 1.2.1 is more effective to increase the ‘sparsity’ of the solution, meaning that it will have less nonzero components.)

(i) Show that (use Exercise 1.2.2), denoting  $\text{sgn}(\mathbf{H}) = \mathbf{1}(\mathbf{H} \geq 0) - \mathbf{1}(\mathbf{H} < 0)$ ,

$$\frac{\partial}{\partial \mathbf{H}} \ell = 2(\mathbf{W}^T \mathbf{W} \mathbf{H} - \mathbf{W}^T \mathbf{X}) + \lambda_1 \text{sgn}(\mathbf{H}) + 2\lambda_2 \mathbf{H}, \quad \frac{\partial^2}{\partial \mathbf{H}^2} \ell = 2(\mathbf{W}^T \mathbf{W} + \lambda_2 \mathbf{I}). \quad (246)$$

(ii) From (i), conclude that the quadratic function  $\mathbf{H} \mapsto \ell(\mathbf{H})$  is convex, and hence its every critical point is a local minimum. It follows that its every critical over the convex constraint set  $\mathbb{R}_{\geq 0}^{r \times n}$  is a global minimum (but not necessarily unique). (See Ref)

**Exercise 4.2.3.** Give a pseudocode of a Block Coordinate Descent algorithm for the following matrix factorization problem

$$\inf_{\mathbf{W} \in \mathbb{R}_{\geq 0}^{d \times r}, \mathbf{H} \in \mathbb{R}_{\geq 0}^{r \times n}} \|\mathbf{X} - \mathbf{WH}\|_F^2 + a_0 \|\mathbf{H}\|_1 + a_1 \|\mathbf{W}\|_1 + a_2 \|\mathbf{W}\|_F \quad (247)$$

The algorithm should be similar to the Alternating Least Squares (ALS) in Algorithm 2, but the subproblems of finding  $\mathbf{H}_n$  and  $\mathbf{W}_n$  should be more detailed by using gradient descent algorithms. (Remark: You may use the gradient given in Exercise 4.2.2. (Hint: The full algorithm is implemented in [Jupyter Notebook](#)))

### 3. Applications of Matrix Factorization

In this section, we apply matrix factorization (MF) to various image and text datasets for the purpose of dictionary learning. See [Jupyter Notebook](#) for the experiments given in this section.

**3.1. MNIST images.** We first apply MF to MNIST dataset to learn 25 dictionary images. There are total four ways to put nonnegativity constraints to  $\mathbf{W}$  and  $\mathbf{H}$ . In the unconstrained case where  $\mathbf{W}$  and  $\mathbf{H}$  have no constraints, the basis images (see Figure 47 left) themselves can take negative values and they are supposed to span the original MNIST images by using any linear combination. Next, if we give nonnegativity constraint on the code  $\mathbf{H}$ , then the basis images (see Figure 47 middle) are still allowed to take negative pixel values but they can only additively combined to span original images. In this case, pixel values in dictionary images can still cancel each other when spanning the original image.

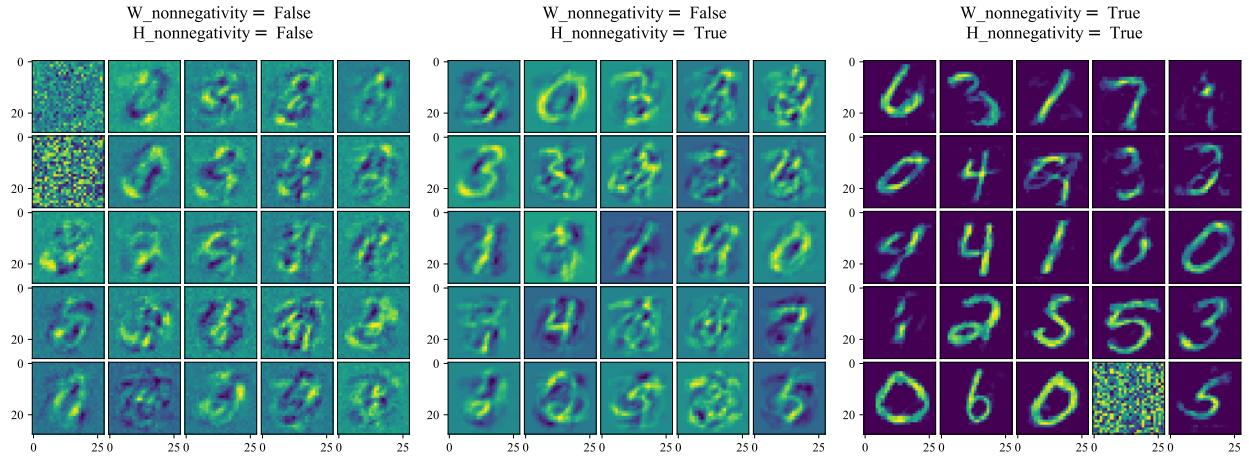


FIGURE 47. Examples of 25 dictionary images learned from MNIST dataset by Matrix Factorization with various nonnegativity constraints.

Lastly, if we constrain both  $\mathbf{H}$  and  $\mathbf{W}$  to be nonnegative, then there is absolutely no cancellation in any linear approximation of the original images by the basis images. Hence the basis images in this case

(see Figure 47) can only have positive pixel values that really contribute to approximating the original images. In other words, the features we see in the dictionary images are true features of the original images. This is the interpretability of NMF that made it so popular in the machine learning and signal processing community [LS99].

In Figure 48, we show similar experiments on MNIST with random padding of thickness 20. We see similar patterns in dictionary images depending on the constraints as in Figure 47.

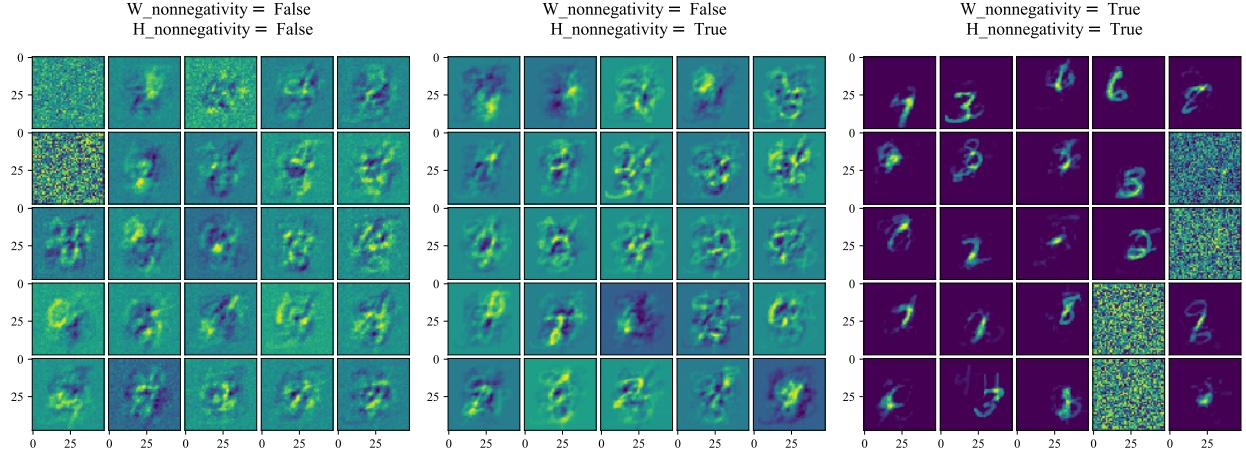


FIGURE 48. Examples of 25 dictionary images learned from MNIST dataset by Matrix Factorization with various nonnegativity constraints. Original  $28 \times 28$  images are padded with zeros randomly so that they are  $48 \times 48$  images.

**3.2. Olivetti face images.** The [Olivetti face dataset](#) contains a set of face images taken between April 1992 and April 1994 at AT&T Laboratories Cambridge. There are ten different images of each of 40 distinct subjects. For some subjects, the images were taken at different times, varying the lighting, facial expressions (open / closed eyes, smiling / not smiling) and facial details (glasses / no glasses). All the images were taken against a dark homogeneous background with the subjects in an upright, frontal position (with tolerance for some side movement) [SH94]. See Figure 49 for some sample images from this dataset.

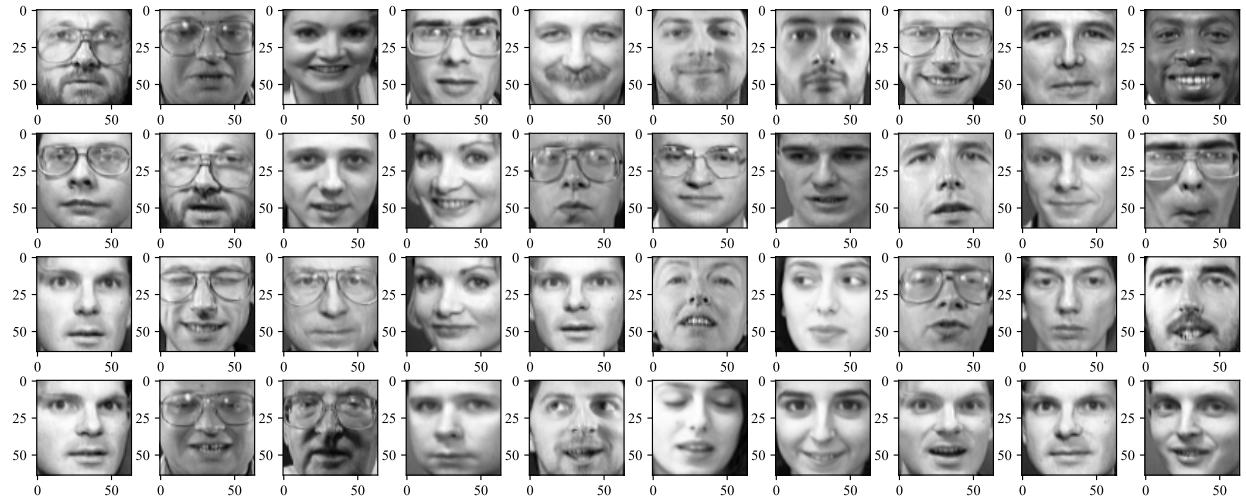


FIGURE 49. Sample images from the [Olivetti face dataset](#)

We perform a similar dictionary learning experiment for the Olivetti face dataset in Figure 50. We see similar patterns in dictionary images depending on the constraints as in Figure 47. Note that the “basis faces” learned by NMF (Figure 50 right) are not only the most interpretable, but also have the feature of *localization*. This ability of being able to learn localized features was first reported in [LS99] in 1999 by Lee and Seung. Later in 2004, Hoyer [Hoy04] further investigation on promoting such localized dictionary by using a particular sparseness regularizer.

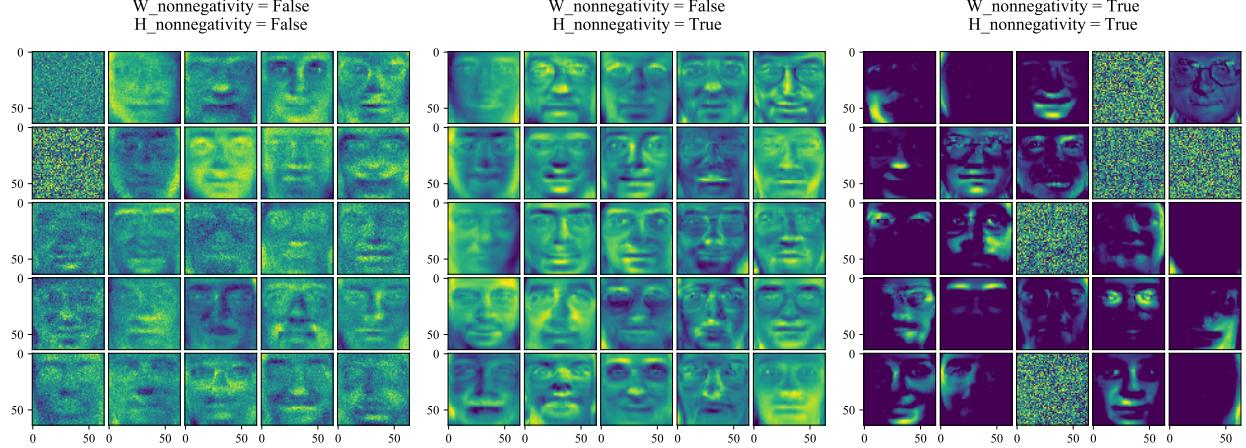


FIGURE 50. Examples of 25 dictionary images learned from the Olivetti face dataset by Matrix Factorization with various nonnegativity constraints.

Below in Figure 51, we discuss the effect of using regularizers on the dictionary matrix  $\mathbf{W}$  in NMF. We use basic  $L_1 + L_2$  regularizer for  $\mathbf{W}$ , which means that we are solving the following optimization problem in (247). The ALS algorithm for solving (247) is similar to Algorithm 2. In fact, the update for  $\mathbf{H}$  is the same, and for updating  $\mathbf{W}$ , we need to solve the following coding problem:

$$\mathbf{W}_n \leftarrow \underset{\mathbf{W} \in \mathbb{R}_{\geq 0}^{d \times r}}{\operatorname{argmin}} \|\mathbf{X} - \mathbf{W}\mathbf{H}_n\|_F^2 + a_1 \|\mathbf{W}\|_1 + a_2 \|\mathbf{W}\|_F^2. \quad (248)$$

Note that (248) is still a (constrained) convex optimization problem that can be solved iteratively by, e.g., gradient descent<sup>4</sup>.

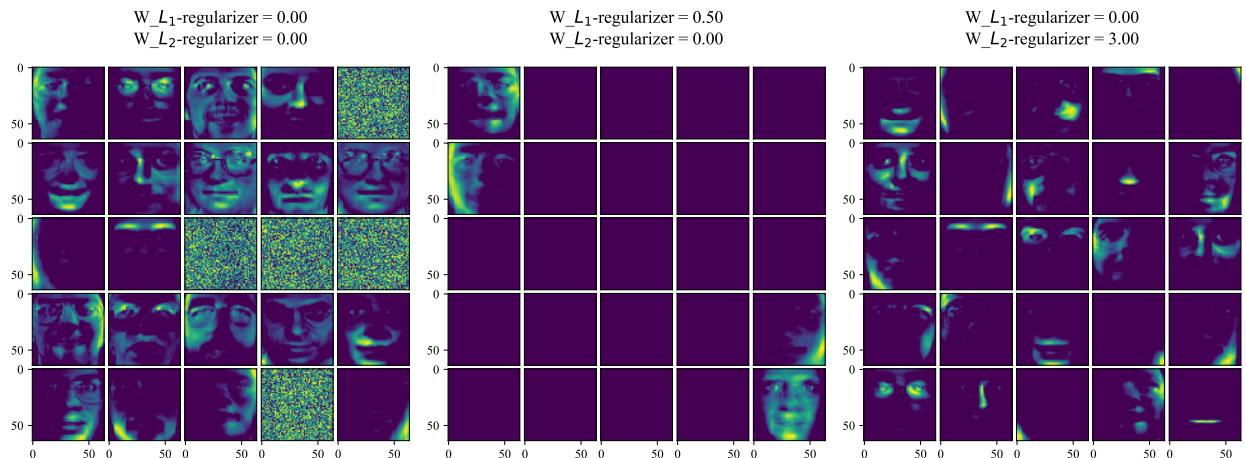


FIGURE 51. Examples of 25 dictionary images learned from the Olivetti face dataset by Matrix Factorization with various regularizers for the dictionary matrix.

<sup>4</sup>It is often the case that one does not have to spend too much iterations to solve sub-problems. In our implementation of ALS for MF in [Jupyter Notebook](#), we used only up to 5 gradient descent iterations for the sub-problems of coding  $\mathbf{H}_n$  and  $\mathbf{W}_n$ .

In Figure 51, observe that using  $L_1$ -regularizer on  $\mathbf{W}$  (see Figure 51 middle) makes most of its columns (i.e., vectorizations of basis images) to be zero, and its nonzero columns are not necessarily localized. On the other hand, using  $L_2$ -regularizer as in (Figure 51 right), we see the columns of  $\mathbf{W}$  are nonzero but they are all localized. One often prefers to learn nonzero localized dictionary elements so that we get more simple and interpretable basis features; On the other hand, one often prefers to have the code matrix  $\mathbf{H}$  to have a few nonzero coordinates so that we only need to use a few dictionary elements to approximate original vectors (e.g., for image reconstruction and denoising). The experiment in Figure 51 suggests that for this, we may use  $L_1$  regularizer for the code  $\mathbf{H}$  and  $L_2$ -regularizer for the dictionary  $\mathbf{W}$ .

**Exercise 4.3.1.** Give an explanation on the effect of  $L_1$ - and  $L_2$ -regularization on the dictionary matrix  $\mathbf{W}$  as we see in Figure 51. Drawing contour plots of  $L_1$ - and  $L_2$ -norm of 2-dimensional vectors would be helpful. (Ref. [Bis06, Sec. 3.1.4].)

**Exercise 4.3.2.** Use the code in Jupyter Notebook that generates Figures 50 and 51 to determine the *largest  $r \in \mathbb{N}$  such that every dictionary image properly converges to some non-random image*, where the nonnegativity constraints and regularizers are chosen as in the six cases in Figures 50 and 51 (estimate such  $r$  in each of the six cases.)

**3.3. Topic modeling for 20Newsgroups.** Next, we apply NMF to the 20Newsgroups dataset (see Section 8.1). Recall that each article is represented as a feature vector of word frequency using a fixed vocabulary (sklearn's standard vocab of size 45534) either by using bag-of-words or tf-idf vectorizer. We use 4483 articles from 10 categories (see (165)) combined. The input data matrix  $\mathbf{X}$  has size  $45534 \times 4483$  and we use NMF to factorize this as  $\mathbf{WH}$ , where  $\mathbf{W} \in \mathbb{R}_{\geq 0}^{45534 \times 10}$  and  $\mathbf{H} \in \mathbb{R}_{\geq 0}^{10 \times 45534}$ .

In case of factorizing data matrix of vectorized images, the columns of  $\mathbf{W}$  were dictionary images that form a basis for the images in the input. For factorizing vectorized texts here, the columns of  $\mathbf{W}$  still form a basis for the texts in the input, but they can also be thought of as *topics* in the entire text corpus. Namely, they give groups of words that are used to (approximately) compose each article nonnegatively.



FIGURE 52. 10 topics learned from the 20Newsgroups data by NMF shown as wordclouds. Top 7 words in each topic are shown.

In Figure 52, we show ten topics learned by NMF as wordclouds. We can indeed see that in each topic, words that belong to the same category appear together (e.g., ‘jesus’, ‘bible’, ‘god’ in topic 6 and ‘os’, ‘pc’, ‘bit’, ‘mac’ in topic 7). It is important to note, however, that we never used label information and these topics are learned in a purely unsupervised manner. NMF is one of the most widely used technique in topic modeling for text data<sup>5</sup>.

<sup>5</sup>Other techniques include Latent Dirichlet Analysis and Nonnegative CP Tensor Decomposition.

We remark that using nonnegativity constraint both on the dictionary  $\mathbf{W}$  and the code  $\mathbf{H}$  are crucial. First, the columns of  $\mathbf{W}$  needs to be vectors of nonnegative entries in order to be interpreted as word frequency vectors. Also, we need to prevent cancellation in the factorization by also making  $\mathbf{H}$  nonnegative so that the words in each topic do share similar context. Indeed, Figure 53 shows the topics learned by a matrix factorization with nonnegative dictionary but unconstrained code. We can see that the words in each topic does not always share the same context when  $\mathbf{H}$  is unconstrained.



FIGURE 53. 10 topics learned from the 20Newsgroups data by a matrix factorization with non-negative dictionary and unconstrained code. Top 7 words in each topic are shown in wordclouds.

**Exercise 4.3.3.** The topics learned by NMF in Figures 52 and 53 uses tf-idf vectorizer of the documents. Use bag-of-words vectorizer instead and reproduce these figures. Can you see any difference? Use the code in [Jupyter Notebook](#) that generates these plots.

#### 4. Principal Component Analysis

In this section, we discuss another matrix-factorization-based classical dimension reduction technique called the ‘principal component analysis’. Roughly speaking, the idea is to project  $n$  observed  $d$ -dimensional data vectors onto a lower  $r$ -dimensional space in a ‘nice way’. We will discuss two equivalent such formulations of ‘variance maximization’ and ‘error minimization’.

**4.1. Maximum variance formulation of PCA.** An important concept in the development of PCA is projection of a vector onto a linear subspace generated by a set of vectors. The details are provided in the following exercise.

**Exercise 4.4.1** (Projection). Let  $\mathbf{x} \in \mathbb{R}^d$  and  $C \subseteq \mathbb{R}^d$ . The projection  $\text{Proj}_C(\mathbf{x})$  of  $\mathbf{x}$  onto  $C$  is defined by

$$\text{Proj}_C(\mathbf{x}) := \underset{\mathbf{u} \in C \subseteq \mathbb{R}^d}{\operatorname{argmin}} \|\mathbf{x} - \mathbf{u}\|_F^2. \quad (249)$$

Namely,  $\text{Proj}_C(\mathbf{x})$  is the point in  $C$  that is closest to  $\mathbf{x}$ . It is defined as a constrained quadratic minimization as above.

(i) Write  $\|\mathbf{x} - \mathbf{u}\|_F^2 = \text{tr}((\mathbf{x} - \mathbf{u})^T (\mathbf{x} - \mathbf{u})) = \text{tr}(\mathbf{u}^T \mathbf{u}) - 2\text{tr}(\mathbf{x}^T \mathbf{u}) + \text{tr}(\mathbf{x}^T \mathbf{x})$ . Show that

$$\frac{\partial}{\partial \mathbf{u}} \|\mathbf{x} - \mathbf{u}\|_F^2 = 2(\mathbf{u} - \mathbf{x}). \quad (250)$$

(ii) Assume  $C$  is the one-dimensional span  $\langle \mathbf{v} \rangle := \{a\mathbf{v} \mid a \in \mathbb{R}\}$ . Then

$$\text{Proj}_{\langle \mathbf{v} \rangle}(\mathbf{x}) := \text{Proj}_{\langle \mathbf{v} \rangle}(\mathbf{x}) = \arg \min_{a \in \mathbb{R}} \|\mathbf{x} - a\mathbf{v}\|_F^2. \quad (251)$$

Use (i) and Lagrange multiplier or solve the above unconstrained optimization directly and show

$$\text{Proj}_{\langle \mathbf{v} \rangle}(\mathbf{x}) = (\mathbf{x}^T \mathbf{v}) \mathbf{v}. \quad (252)$$

(iii) Let  $\mathbf{W} = [\mathbf{u}_1, \dots, \mathbf{u}_r] \in \mathbb{R}^{d \times r}$  be such that its columns are orthonormal, i.e.,  $\mathbf{W}^T \mathbf{W} = \mathbf{I} \in \mathbb{R}^{r \times r}$ . Show that

$$\text{Proj}_{\langle \mathbf{W} \rangle}(\mathbf{x}) := \text{Proj}_{\langle \mathbf{u}_1, \dots, \mathbf{u}_r \rangle}(\mathbf{x}) = \sum_{i=1}^r (\mathbf{u}_i^T \mathbf{x}) \mathbf{u}_i = \mathbf{W}(\mathbf{W}^T \mathbf{x}). \quad (253)$$

Next, we introduce the following notions:

**Mean and Variance:** Fix  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$ . We define its *mean* and *variance* as

$$\bar{\mathbf{X}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \in \mathbb{R}^d, \quad \text{Var}(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \bar{\mathbf{X}}\|_F^2 \in \mathbb{R}. \quad (254)$$

**Covariance:** The *covariance matrix*  $\mathbf{S} = \mathbf{S}(\mathbf{X}) \in \mathbb{R}^{d \times d}$  of  $\mathbf{X} \in \mathbb{R}^{d \times n}$  is defined by

$$\mathbf{S} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{X}})(\mathbf{x}_i - \bar{\mathbf{X}})^T. \quad (255)$$

**Projection:** Fix  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$  and a dictionary matrix  $\mathbf{W} \in \mathbb{R}^{d \times r}$  with orthonormal columns (i.e.,  $\mathbf{W}^T \mathbf{W} = \mathbf{I} \in \mathbb{R}^{r \times r}$ ). Define the *projection of  $\mathbf{X}$  onto the subspace spanned by  $\mathbf{W}$*  as

$$\text{Proj}_{\langle \mathbf{W} \rangle}(\mathbf{X}) := [\mathbf{W} \mathbf{W}^T \mathbf{x}_1, \dots, \mathbf{W} \mathbf{W}^T \mathbf{x}_n] = \mathbf{W}(\mathbf{W}^T \mathbf{X}) \in \mathbb{R}^{d \times n}. \quad (256)$$

According to Exercise 4.4.1, the columns of  $\text{Proj}_{\langle \mathbf{W} \rangle}(\mathbf{X})$  agree with the projection of the columns of  $\mathbf{X}$  onto  $\langle \mathbf{W} \rangle$ .<sup>6</sup>

Projection inherently loses information by collapsing two distinct data points closer to each other. Hence if we were to find the best  $r$ -dimensional projection of the data matrix  $\mathbf{X}$  that preserves as much information as possible, we should keep the variance of the  $n$  vectors after projection as large as possible. Moreover, looking for  $r$ -dimensional basis to project onto, we will require the basis vectors to have unit norm. Hence we would like to seek a matrix  $\mathbf{W} \in \mathbb{R}^{d \times r}$  of  $r$  basis vectors of unit norm such that the variance after projection is maximized:

$$\begin{aligned} & \text{(PCA; maximum variance)} \quad \hat{\mathbf{W}} = \arg \max_{\substack{\mathbf{W} \in \mathbb{R}^{d \times r} \\ \text{rank}(\mathbf{W})=r}} \text{Var}(\text{Proj}_{\langle \mathbf{W} \rangle}(\mathbf{X})), \end{aligned} \quad (257)$$

where  $\text{rank}(\mathbf{W})$  denotes the dimension of the linear subspace  $\langle \mathbf{W} \rangle$  generated by the columns of  $\mathbf{W}$ .<sup>7</sup>

In the following proposition, we show that finding the top  $r$  eigenvectors of the covariance matrix  $\mathbf{S}$  solves the PCA problem in (257)<sup>8</sup>. These vectors are also called the *principal components* of  $\mathbf{X}$ , and the associated eigenvalues are called *singular values*.

**Proposition 4.4.2.** Let  $\mathbf{X} \in \mathbb{R}^{d \times n}$  and  $\mathbf{S} = \mathbf{S}(\mathbf{X}) \in \mathbb{R}^{d \times d}$  denote its covariance matrix. Fix  $r \in \{1, \dots, n\}$ . Then a solution to the maximum variance PCA problem in (262) is given by  $\mathbf{W} \in \mathbb{R}^{d \times r}$  where

$$\text{The columns of } \mathbf{W} \text{ are the unit eigenvectors of } \mathbf{S} \text{ with top } r \text{ largest eigenvalues.} \quad (258)$$

Furthermore, the maximum variance  $\text{Var}(\mathbf{W}^T \mathbf{X})$  after projection achieved equals the sum of the top  $r$  largest eigenvalues of  $\mathbf{S}$ .

<sup>6</sup>Note that  $\mathbf{W}^T \mathbf{X} \in \mathbb{R}^{r \times d}$  gives the  $r$ -dimensional compressed version of  $\mathbf{X}$ , whose columns give coordinates of the columns of  $\mathbf{X}$  after projection w.r.t. the orthonormal basis of the columns of  $\mathbf{W}$ .

<sup>7</sup>Note that  $\text{rank}(\mathbf{W}) = \dim(\langle \mathbf{W} \rangle)$ . So the columns of  $\mathbf{W}$  are linearly independent iff  $\text{rank}(\mathbf{W}) = r$

<sup>8</sup>But there could be other equivalent solutions.

PROOF. First, write  $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_r]$  and we will choose an orthonormal basis  $\mathbf{u}_1, \dots, \mathbf{u}_r$  for the linear subspace  $\langle \mathbf{w}_1, \dots, \mathbf{w}_r \rangle$  spanned by the columns of  $\mathbf{W}$  (e.g., by Gram-Schmidt). Then since  $\langle \mathbf{w}_1, \dots, \mathbf{w}_r \rangle = \langle \mathbf{u}_1, \dots, \mathbf{u}_r \rangle$ , we have

$$\text{Proj}_{\langle \mathbf{W} \rangle}(\mathbf{X}) = \text{Proj}_{\langle \mathbf{u}_1, \dots, \mathbf{u}_r \rangle}(\mathbf{X}). \quad (259)$$

Hence without loss of generality, we may assume that the columns of  $\mathbf{W}$  are orthonormal. i.e.,  $\mathbf{W}^T \mathbf{W} = \mathbf{I} \in \mathbb{R}^{r \times r}$ , in (257).

Since the columns of  $\mathbf{W}$  are orthonormal, note that

$$\text{Var}(\text{Proj}_{\langle \mathbf{W} \rangle}(\mathbf{X})) = \text{Var}(\mathbf{W}\mathbf{W}^T \mathbf{X}) = \frac{1}{n} \sum_{i=1}^n \|\mathbf{W}\mathbf{W}^T(\mathbf{x}_i - \bar{\mathbf{X}})\|_F^2 = \text{tr}(\mathbf{W}\mathbf{W}^T \mathbf{S} \mathbf{W}\mathbf{W}^T), \quad (260)$$

where  $\mathbf{S} = \mathbf{S}(\mathbf{X}) \in \mathbb{R}^{d \times d}$  denotes the covariance matrix of  $\mathbf{X}$  defined at (255). By the cyclic property of trace (see Exercise 4.4.3) and since  $\mathbf{W}^T \mathbf{W} = \mathbf{I}$ ,

$$\text{tr}(\mathbf{W}\mathbf{W}^T \mathbf{S} \mathbf{W}\mathbf{W}^T) = \text{tr}(\mathbf{W}^T \mathbf{S} \mathbf{W}(\mathbf{W}^T \mathbf{W})) = \text{tr}(\mathbf{W}^T \mathbf{S} \mathbf{W}). \quad (261)$$

Thus (257) is equivalent to

$$\hat{\mathbf{W}} = \underset{\substack{\mathbf{W} \in \mathbb{R}^{d \times r} \\ \mathbf{W}^T \mathbf{W} = \mathbf{I}}}{\arg \max} \text{tr}(\mathbf{W}^T \mathbf{S} \mathbf{W}). \quad (262)$$

We claim the following:<sup>9</sup>

If  $\mathbf{W}$  is a stationary point of (262) then we have an orthonormal  $\mathbf{U} \in \mathbb{R}^{r \times r}$  s.t.

- (1)  $\mathbf{U}^T \mathbf{W} \in \mathbb{R}^{d \times r}$  is orthonormal; and
- (2)  $\text{tr}(\mathbf{W}^T \mathbf{S} \mathbf{W}) = \text{tr}(\mathbf{U} \mathbf{W}^T \mathbf{S} \mathbf{W} \mathbf{U}^T)$ ; and
- (3)  $\mathbf{S}(\mathbf{W} \mathbf{U}^T) = (\mathbf{W} \mathbf{U}^T) \mathbf{D}$  where  $\mathbf{D} \in \mathbb{R}^{r \times r}$  is diagonal with  $r$  eigenvalues of  $\mathbf{S}$  for its diagonal entries.

In particular, this claim implies that if  $\hat{\mathbf{W}} \in \mathbb{R}^{d \times r}$  is a solution to (262), then we can find an ‘equivalent maximizer’  $\mathbf{U}^T \hat{\mathbf{W}}$  (by (1) and (2)) with the additional property that the columns of  $\mathbf{U}^T \hat{\mathbf{W}}$  are eigenvectors of  $\mathbf{S}$  (by (3)). Notice that, by writing  $\mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_r)$  where  $\lambda_i$ ’s are eigenvalues of  $\mathbf{S}$ ,

$$\text{tr}((\mathbf{W} \mathbf{U}^T)^T \mathbf{S}(\mathbf{W} \mathbf{U}^T)) = \text{tr}((\mathbf{W} \mathbf{U}^T)^T (\mathbf{W} \mathbf{U}^T) \mathbf{D}) \quad (263)$$

$$= \text{tr}(\mathbf{U} \mathbf{W}^T \mathbf{W} \mathbf{U}^T \mathbf{D}) \quad (264)$$

$$= \text{tr}(\mathbf{D}) = \lambda_1 + \dots + \lambda_r. \quad (265)$$

Therefore, if we choose the columns of  $\mathbf{U}^T \mathbf{W}$  to be an orthonormal set of  $r$  eigenvectors of  $\mathbf{S}$  with top  $r$  eigenvalues of  $\mathbf{S}$ , then  $\mathbf{U}^T \mathbf{W}$  has to be a solution to (262), hence proving the assertion.

It remains to show the claim above. Two essential ingredients are multivariate Lagrange multipliers and spectral theorem for real symmetric matrices.

In order to find all stationary points of (262), we formulate its Lagrangian as

$$L(\mathbf{W}; \boldsymbol{\lambda}) := \text{tr}(\mathbf{W}^T \mathbf{S} \mathbf{W}) + \sum_{j=1}^r \lambda_{jj} (1 - \text{tr}(\mathbf{W}[:, j]^T \mathbf{W}[:, j])) + \sum_{\substack{1 \leq i, j \leq r \\ i \neq j}} \lambda_{ij} (0 - \text{tr}(\mathbf{W}[:, i]^T \mathbf{W}[:, j])), \quad (266)$$

where  $\boldsymbol{\lambda} = (\lambda_{ij}) \in \mathbb{R}^{r \times r}$  is a dual variable<sup>10</sup>. By the method of multivariate Lagrange multiplier, a matrix  $\mathbf{W} \in \mathbb{R}^{d \times r}$  is a stationary point of (262) only if there exists  $\boldsymbol{\lambda} \in \mathbb{R}^{r \times r}$  such that it is a stationary point of  $L(\cdot; \boldsymbol{\lambda})$ . So, differentiating the objective function above by the  $j$ th column  $\mathbf{W}[:, j]$  of  $\mathbf{W}$ ,

$$\frac{\partial}{\partial \mathbf{W}[:, j]} = \mathbf{S} \mathbf{W}[:, j] - \sum_{1 \leq i \leq r} \lambda_{ij} \mathbf{W}[:, i] \quad (267)$$

$$= \mathbf{S} \mathbf{W}[:, j] - \mathbf{W}[:, i] \boldsymbol{\lambda}[:, j]. \quad (268)$$

---

<sup>9</sup>Simply speaking, it says that there exists a global maximizer  $\hat{\mathbf{W}}_*$  of (262) such that its columns are eigenvectors of  $\mathbf{S}$ .

<sup>10</sup>Just think of it as a multi-dimensional Lagrange multiplier

Hence, in order for  $\mathbf{W} \in \mathbb{R}^{d \times r}$  to be a stationary point of (262), there should exist a dual variable  $\boldsymbol{\lambda} \in \mathbb{R}^{r \times r}$  such that

$$\mathbf{S}\mathbf{W} = \mathbf{W}\boldsymbol{\lambda}. \quad (269)$$

Using  $\mathbf{W}^T\mathbf{W} = \mathbf{I} \in \mathbb{R}^{r \times r}$ , this implies

$$\mathbf{W}^T\mathbf{S}\mathbf{W} = \boldsymbol{\lambda}. \quad (270)$$

Next, notice that  $\mathbf{W}^T\mathbf{S}\mathbf{W} = \boldsymbol{\lambda}$  is a real symmetric matrix. So by a spectral theorem (see Exercise 4.4.4), there exists an orthonormal  $\mathbf{U} \in \mathbb{R}^{r \times r}$ <sup>11</sup> such that  $\mathbf{W}^T\mathbf{S}\mathbf{W} = \mathbf{U}^T\mathbf{D}\mathbf{U}$ , where  $\mathbf{D}$  is a diagonal matrix whose diagonal entries are the eigenvalues of  $\mathbf{W}^T\mathbf{S}\mathbf{W} \in \mathbb{R}^{r \times r}$ . This shows

$$\mathbf{S}(\mathbf{W}\mathbf{U}^T) = (\mathbf{W}\mathbf{U}^T)\mathbf{D}, \quad (\text{claim (2)}). \quad (271)$$

Also, it is easy to see that  $\mathbf{U}^T\mathbf{W} \in \mathbb{R}^{d \times r}$  is orthonormal, since  $(\mathbf{U}^T\mathbf{W})^T(\mathbf{U}^T\mathbf{W}) = \mathbf{W}^T\mathbf{U}\mathbf{U}^T\mathbf{W} = \mathbf{W}^T\mathbf{W} = \mathbf{I} \in \mathbb{R}^{r \times r}$  (claim (1)). Lastly, observe that using  $\mathbf{U}^T\mathbf{U} = \mathbf{I}$ ,

$$\text{tr}(\mathbf{W}^T\mathbf{S}\mathbf{W}) = \text{tr}(\mathbf{W}^T\mathbf{S}\mathbf{W}\mathbf{U}\mathbf{U}^T) = \text{tr}(\mathbf{U}\mathbf{W}^T\mathbf{S}\mathbf{W}\mathbf{U}^T) \quad (\text{claim (3)}). \quad (272)$$

This shows the claim, thereby finishing the proof of the assertion.  $\square$

**Exercise 4.4.3** (Cyclic property of trace). Show that  $\text{tr}(\mathbf{ABC}) = \text{tr}(\mathbf{BCA})$ .

The  $r$  principal components of  $\mathbf{X}$  as stated in Proposition 4.4.2 are always orthogonal. This is implied by the spectral theorem for real symmetric matrices, see Exercise 4.4.4.

**Exercise 4.4.4** (Spectral Theorem for Real Symmetric Matrices). Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  be a real symmetric matrix, i.e.,  $\mathbf{A}^T = \mathbf{A}$ . Then the spectral theorem for real symmetric matrices states that the following hold:

1.  $\mathbf{A}$  has  $n$  real eigenvalues (counting multiplicities);
2. There exists an orthonormal matrix  $\mathbf{U} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{U}^T\mathbf{U} = \mathbf{I}$  and a diagonal matrix  $\mathbf{D} \in \mathbb{R}^{n \times n}$  such that

$$\mathbf{A} = \mathbf{U}^T\mathbf{D}\mathbf{U}; \quad (273)$$

3. In (ii), the orthonormal matrix  $\mathbf{U}$  consists of eigenvectors of  $\mathbf{A}$  and the diagonal entries of  $\mathbf{D}$  are the eigenvalues of  $\mathbf{A}$ ;

In this exercise, we will prove this result by an induction on  $n$ .

- (i) Show that all of the eigenvalues of  $\mathbf{A}$  are real.
- (ii) Let  $\mathbf{v}_1$  be an eigenvector of  $\mathbf{A}$  and let  $\lambda_1$  be its associated eigenvalue. By (i),  $\lambda_1$  is real. Extend  $\mathbf{v}_1$  to an orthonormal basis  $\mathcal{B} := \{\mathbf{v}_1, \mathbf{u}_2, \dots, \mathbf{u}_n\}$  for  $\mathbb{R}^n$  and let  $\mathbf{Q} := [\mathbf{v}_1, \mathbf{u}_2, \dots, \mathbf{u}_n] \in \mathbb{R}^{n \times n}$ . Show that

$$\mathbf{Q}^T\mathbf{A}\mathbf{Q} = \left[ \begin{array}{c|cccc} \lambda_1 & 0 & \dots & 0 \\ \hline 0 & & & & \\ \vdots & & \mathbf{A}' & & \\ 0 & & & & \end{array} \right], \quad (274)$$

where  $\mathbf{A}' \in \mathbb{R}^{(n-1) \times (n-1)}$  is a real symmetric matrix.

- (iii) Suppose there exists an orthonormal matrix  $\mathbf{V} \in \mathbb{R}^{(n-1) \times (n-1)}$  and a diagonal matrix  $\mathbf{D}' \in \mathbb{R}^{(n-1) \times (n-1)}$  such that  $\mathbf{A}' = \mathbf{V}^T\mathbf{D}'\mathbf{V}$ . Define

$$\bar{\mathbf{V}} = \left[ \begin{array}{c|cccc} 1 & 0 & \dots & 0 \\ \hline 0 & & & & \\ \vdots & & \mathbf{V} & & \\ 0 & & & & \end{array} \right] \in \mathbb{R}^{n \times n}. \quad (275)$$

---

<sup>11</sup> $\mathbf{U}^T\mathbf{U} = \mathbf{U}\mathbf{U}^T = \mathbf{I}$  since  $\mathbf{U}$  is a square matrix.

Then show that

$$\mathbf{A} = \mathbf{Q} \bar{\mathbf{V}}^T \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & & & \\ \vdots & & \mathbf{D}' & \\ 0 & & & \end{bmatrix} \bar{\mathbf{V}} \mathbf{Q}^T. \quad (276)$$

Also show that  $\mathbf{U} := \bar{\mathbf{V}} \mathbf{Q}^T \in \mathbb{R}^{n \times n}$  is orthonormal.

- (iv) By using (i)-(ii), prove that if  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is symmetric, show that there are orthogonal and diagonal matrices  $\mathbf{U}, \mathbf{D} \in \mathbb{R}^{n \times n}$  such that  $\mathbf{A} = \mathbf{U}^T \mathbf{D} \mathbf{U}$  and  $\mathbf{D}$  consists of the eigenvalues of  $\mathbf{A}$  as its diagonal entries. Also deduce that the columns of  $\mathbf{U}$  are the eigenvectors of  $\mathbf{A}$ .

**4.2. Minimum reconstruction error formulation of PCA.** In this subsection, we give an alternative formulation of PCA by minimizing reconstruction error after projecting onto an  $r$ -dimensional orthonormal basis.

As before, let  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$  be a data matrix. We will center it by subtracting its mean  $\bar{\mathbf{X}}$  so that  $\bar{\mathbf{X}} = \mathbf{0}$  and  $\mathbf{S} = n^{-1} \mathbf{X} \mathbf{X}^T$ . Let  $\mathbf{W} \in \mathbb{R}^{d \times r}$  be such that its columns are orthonormal, i.e.,  $\mathbf{W}^T \mathbf{W} = \mathbf{I} \in \mathbb{R}^{r \times r}$ . The ‘reconstruction error’ of approximating  $\mathbf{x}_i$  by the projection  $\text{Proj}_{\langle \mathbf{W} \rangle}(\mathbf{x}_i)$  is  $\|\mathbf{x}_i - \text{Proj}_{\langle \mathbf{W} \rangle}(\mathbf{x}_i)\|_F^2$ . In the minimum reconstruction error formulation of PCA, we would like to choose  $\mathbf{W}$  such that the total reconstruction error of all columns of  $\mathbf{X}$  is minimized:

$$(\text{PCA; min. recons. error}) \quad \hat{\mathbf{W}} = \underset{\substack{\mathbf{W} \in \mathbb{R}^{d \times r} \\ \text{rank}(\mathbf{W})=r}}{\arg \min} \frac{1}{n} \|\mathbf{X} - \text{Proj}_{\langle \mathbf{W} \rangle}(\mathbf{X})\|_F^2. \quad (277)$$

In the following proposition, we show the above alternative formulation of PCA is equivalent to the maximum variance one in (257).

**Proposition 4.4.5.** *The minimum reconstruction error formulation of PCA (277) is equivalent to the maximum variance formulation of PCA (257).*

PROOF. As in the proof of Proposition 4.4.2, we may assume that  $\mathbf{W} \in \mathbb{R}^{d \times r}$  is orthonormal, i.e.,  $\mathbf{W}^T \mathbf{W} = \mathbf{I} \in \mathbb{R}^{r \times r}$ . We claim that for orthonormal  $\mathbf{W}$ ,

$$\frac{1}{n} \|\mathbf{X} - \text{Proj}_{\langle \mathbf{W} \rangle}(\mathbf{X})\|_F^2 = \text{Var}(\mathbf{X}) - \text{tr}(\mathbf{W}^T \mathbf{S} \mathbf{W}), \quad (278)$$

where  $\mathbf{S} \in \mathbb{R}^{d \times d}$  denotes the covariance matrix of  $\mathbf{X} \in \mathbb{R}^{d \times n}$ . Then the assertion follows from (278) since then (277) is equivalent to

$$\underset{\substack{\mathbf{W} \in \mathbb{R}^{d \times r} \\ \mathbf{W}^T \mathbf{W} = \mathbf{I}}}{\arg \min} [\text{Var}(\mathbf{X}) - \text{tr}(\mathbf{W}^T \mathbf{S} \mathbf{W})] = \underset{\substack{\mathbf{W} \in \mathbb{R}^{d \times r} \\ \mathbf{W}^T \mathbf{W} = \mathbf{I}}}{\arg \max} \text{tr}(\mathbf{W}^T \mathbf{S} \mathbf{W}), \quad (279)$$

and the last is an equivalent formulation of the maximum variance PCA (257) as we saw in (262).

It remains to show the identity (278) for orthonormal  $\mathbf{W} \in \mathbb{R}^{d \times r}$ . To this end, we extend  $\mathbf{W}$  to obtain a complete orthonormal basis (ONB) for  $\mathbb{R}^d$ . Namely, we can choose  $\mathbf{W}^\perp \in \mathbb{R}^{d \times (d-r)}$  be such that the columns in their horizontal concatenation  $\bar{\mathbf{W}} := [\mathbf{W}, \mathbf{W}^\perp] \in \mathbb{R}^{d \times d}$  form an orthonormal basis (ONB) for  $\mathbb{R}^d$  (i.e.,  $\bar{\mathbf{W}}^T \bar{\mathbf{W}} = \mathbf{I} \in \mathbb{R}^{d \times r}$ .) Fix a column  $\mathbf{x}_i$  of  $\mathbf{X}$ . Note that since the columns of  $\bar{\mathbf{W}}$  span the full space  $\mathbb{R}^d$ , we have

$$\mathbf{x}_i = \text{Proj}_{\langle \bar{\mathbf{W}} \rangle}(\mathbf{x}_i) = \bar{\mathbf{W}} \bar{\mathbf{W}}^T \mathbf{x}_i. \quad (280)$$

On the other hand, note that  $\bar{\mathbf{W}}[\mathbf{W}, \mathbf{O}]^T = [\mathbf{W}, \mathbf{W}^\perp][\mathbf{W}, \mathbf{O}]^T = \mathbf{W}\mathbf{W}^T$ , so we have

$$\text{Proj}_{\langle \mathbf{W} \rangle}(\mathbf{x}_i) = \mathbf{W}\mathbf{W}^T \mathbf{x}_i = \bar{\mathbf{W}}[\mathbf{W}, \mathbf{O}]^T \mathbf{x}_i. \quad (281)$$

Thus we get

$$\|\mathbf{x}_i - \text{Proj}_{\langle \mathbf{W} \rangle}(\mathbf{x}_i)\|_F^2 = \|\bar{\mathbf{W}} \bar{\mathbf{W}}^T \mathbf{x}_i - \bar{\mathbf{W}}[\mathbf{W}, \mathbf{O}]^T \mathbf{x}_i\|_F^2 \quad (282)$$

$$= \|\bar{\mathbf{W}}([\mathbf{W}, \mathbf{W}^\perp]^T - [\mathbf{W}, \mathbf{O}]^T) \mathbf{x}_i\|_F^2 \quad (283)$$

$$= \|([\mathbf{W}, \mathbf{W}^\perp]^T - [\mathbf{W}, \mathbf{O}]^T) \mathbf{x}_i\|_F^2 \quad (284)$$

$$= \|(\mathbf{W}^\perp)^T \mathbf{x}_i\|_F^2, \quad (285)$$

where the third equality follows from the fact that  $\bar{\mathbf{W}}$  is orthonormal.<sup>12</sup> Also note that  $\|(\mathbf{W}^\perp)^T \mathbf{x}_i\|_F^2 = \|\mathbf{x}_i^T (\mathbf{W}^\perp)\|_F^2 = (\mathbf{W}^\perp) \mathbf{x}_i \mathbf{x}_i^T (\mathbf{W}^\perp)^T$ . Hence

$$\frac{1}{n} \|\mathbf{X} - \text{Proj}_{\langle \mathbf{W} \rangle}(\mathbf{X})\|_F^2 = \frac{1}{n} \sum_{i=1}^n (\mathbf{W}^\perp) \mathbf{x}_i \mathbf{x}_i^T (\mathbf{W}^\perp)^T = \text{tr}((\mathbf{W}^\perp) \mathbf{S} (\mathbf{W}^\perp)^T). \quad (286)$$

Lastly, note that  $\text{tr}(\bar{\mathbf{W}}^T \mathbf{S} \bar{\mathbf{W}}) = \text{tr}(\mathbf{S} \bar{\mathbf{W}} \bar{\mathbf{W}}^T) = \text{tr}(\mathbf{S}) = \text{Var}(\mathbf{X})$ . Hence we get

$$\text{Var}(\mathbf{X}) = \text{tr}(\bar{\mathbf{W}}^T \mathbf{S} \bar{\mathbf{W}}) = \text{tr}([\mathbf{W}, \mathbf{W}^\perp]^T \mathbf{S} [\mathbf{W}, \mathbf{W}^\perp]) \quad (287)$$

$$= \text{tr}(\mathbf{W}^T \mathbf{S} \mathbf{W}) + \text{tr}((\mathbf{W}^\perp)^T \mathbf{S} (\mathbf{W}^\perp)) + \text{tr}(\mathbf{W}^T \mathbf{S} (\mathbf{W}^\perp)) + \text{tr}((\mathbf{W}^\perp)^T \mathbf{S} \mathbf{W}) \quad (288)$$

$$= \text{tr}(\mathbf{W}^T \mathbf{S} \mathbf{W}) + \text{tr}((\mathbf{W}^\perp)^T \mathbf{S} (\mathbf{W}^\perp)) + \text{tr}\left(\mathbf{S} \underbrace{(\mathbf{W}^\perp \mathbf{W}^T)}_{\mathbf{O}}\right) + \text{tr}\left(\mathbf{S} \underbrace{\mathbf{W} (\mathbf{W}^\perp)^T}_{\mathbf{O}}\right) \quad (289)$$

$$= \text{tr}(\mathbf{W}^T \mathbf{S} \mathbf{W}) + \text{tr}((\mathbf{W}^\perp)^T \mathbf{S} (\mathbf{W}^\perp)). \quad (290)$$

Thus, the claimed identity (278) follows from combining the above and (286). This finishes the proof.  $\square$

## 5. Applications of PCA and Matrix Factorization

In this section, we apply PCA to various image and text datasets for the purpose of dictionary learning and compare with matrix factorization (MF). See [Jupyter Notebook](#) for the experiments given in this section. Here we use sklearn's PCA implementation `sklearn.decomposition.PCA`.

**5.1. MNIST.** We first apply PCA to the MNIST dataset to learn 24 dictionary images (principal components) and their singular values. These are shown in Figures 54. In Figure 54, it appears that the principal components with larger singular values captures the overall shape of the digits whereas those with smaller singular values gives finer details of the shape.

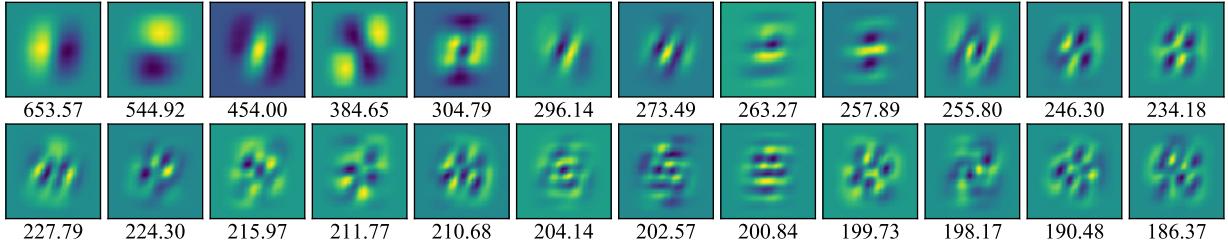


FIGURE 54. Examples of 24 dictionary images (principal components) learned from MNIST dataset by PCA. The numbers below each image are the corresponding singular values.

In Figure 55, we compare the MNIST image dictionary learned by PCA against those learned by MF. The numbers underneath the dictionary images for MF cases represent their ‘total usage’ computed by row sums of the code matrix  $\mathbf{H}$ . Namely, observe that from the factorization  $\mathbf{X} \approx \mathbf{WH} \in \mathbb{R}^{d \times n}$ ,  $\mathbf{H}[i, j]$  equals the linear coefficient of  $\mathbf{W}[:, i]$  in approximating the  $i$ th data point  $\mathbf{X}[:, j]$ . Hence we can define

$$\text{total usage of } \mathbf{W}[:, i] := \sum_{j=1}^n |\mathbf{H}[i, j]|. \quad (291)$$

<sup>12</sup>Suppose  $\mathbf{A} \in \mathbb{R}^{d \times d}$  is orthonormal, i.e.,  $\mathbf{A}^T \mathbf{A} = \mathbf{I}$ . Let  $\mathbf{x} \in \mathbb{R}^d$ . Then  $\|\mathbf{Ax}\|_F^2 = (\mathbf{Ax})^T (\mathbf{Ax}) = \mathbf{x}^T \mathbf{A}^T \mathbf{Ax} = \mathbf{x}^T \mathbf{Ix} = \mathbf{x}^T \mathbf{x} = \|\mathbf{x}\|_F^2$ .

Note that we take the absolute value of the code coefficient  $\mathbf{H}[i, j]$  above in order to prevent cancellation in computing the total usage. Figure 56 shows similar experiments on MNIST with random padding of thickness 20. It is clear that NMF-learned dictionary are the easiest to interpret.

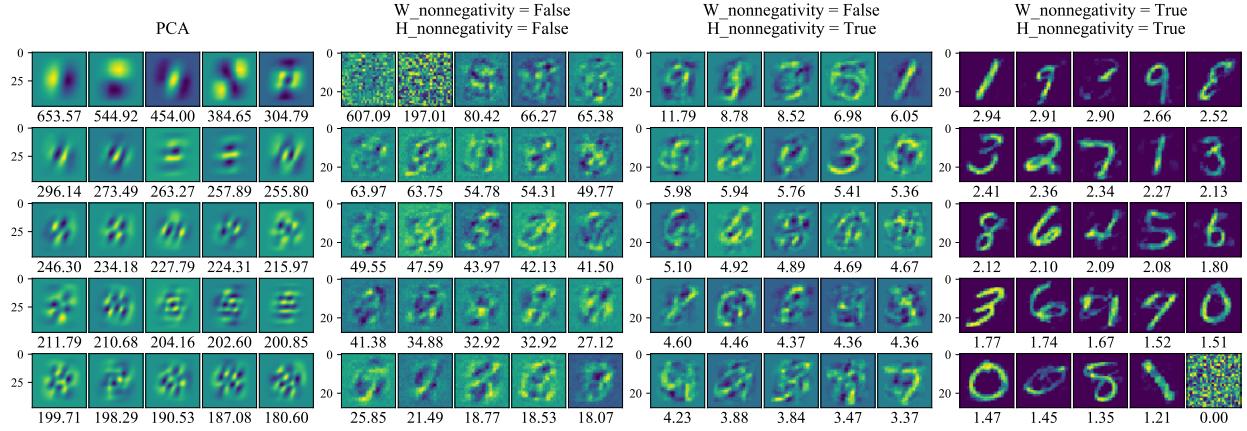


FIGURE 55. Examples of 25 dictionary images learned from MNIST dataset by Matrix Factorization with various nonnegativity constraints.

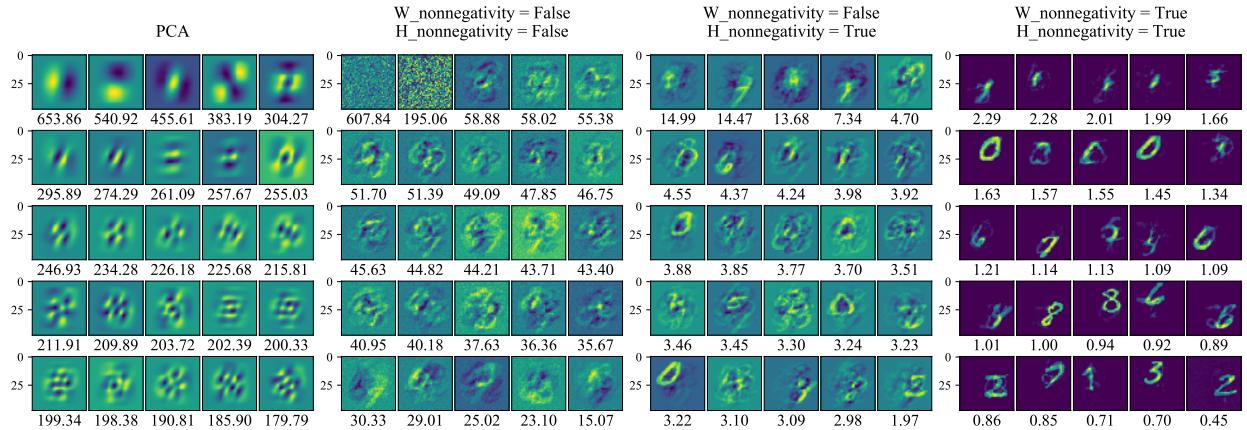


FIGURE 56. Examples of 25 dictionary images learned from MNIST dataset by Matrix Factorization with various nonnegativity constraints. Original  $28 \times 28$  images are padded with zeros randomly so that they are  $48 \times 48$  images.

**5.2. Olivetti face images.** We perform similar dictionary learning experiments for the Olivetti face dataset in Figure 57. Since the principal components may have negative entries, the PCA-learned dictionary face images (a.k.a. ‘eigenfaces’) may have some shapes that are not part of any of the actual face images.

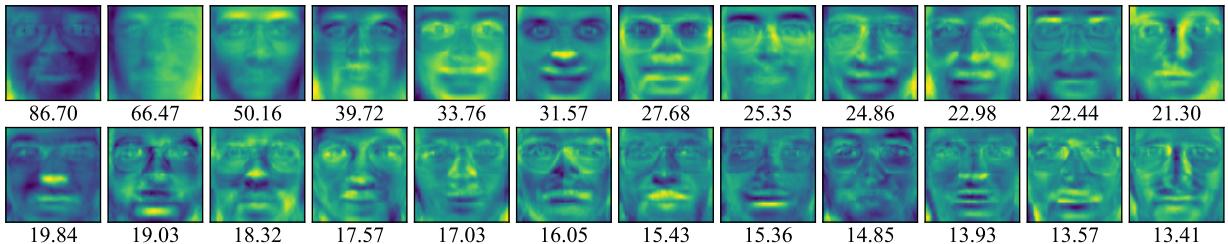


FIGURE 57. Examples of 24 dictionary images (principal components) learned from the Olivetti face dataset by PCA. The numbers below each image are the corresponding singular values.

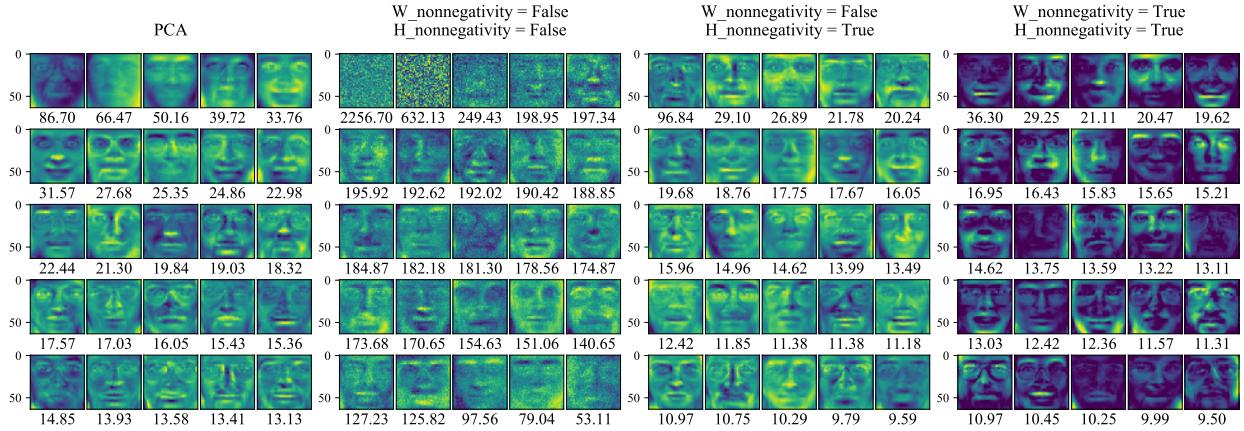


FIGURE 58. Examples of 25 dictionary images learned from the Olivetti face dataset by PCA and Matrix Factorization with various nonnegativity constraints.

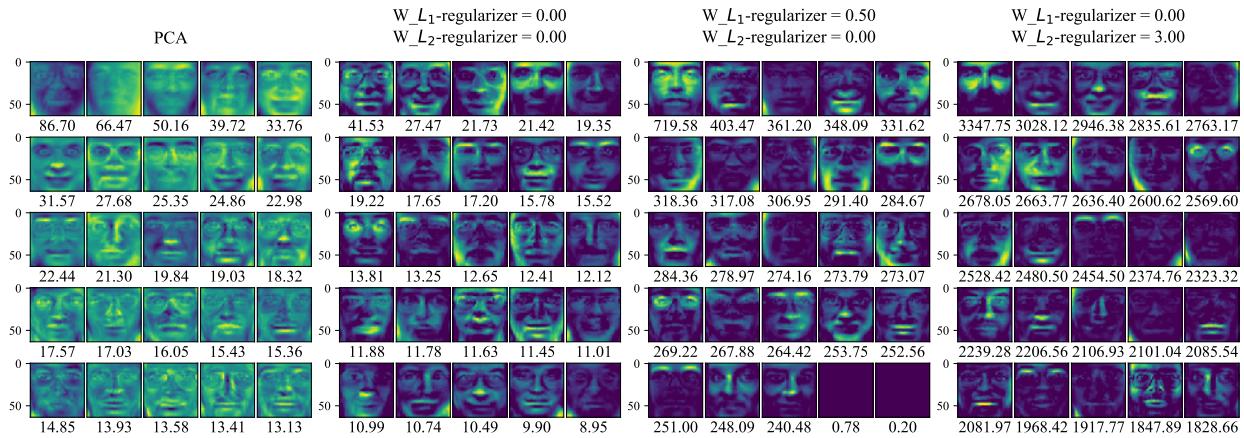


FIGURE 59. Examples of 25 dictionary images learned from the Olivetti face dataset by PCA and Matrix Factorization with various dictionary regularization.

## 6. Probabilistic PCA and EM algorithm

In this section, we will consider a probabilistic version of PCA and an iterative algorithm for solving the PCA optimization problem called the Expectation-Maximization (EM) algorithm. The discussion in this section is based on [Bis06, Sec. 9.3 and 12.2] and [Row98, TB99].

**6.1. The EM algorithm for PCA.** The goal of this section is to derive the following *Expectation-Maximization* (EM) algorithm for PCA. In order to provide the full background and discussion, we need to introduce a probabilistic version of PCA and formulate a maximum likelihood formulation of PCA. The EM algorithms form a general class of algorithms used to solve a complex maximum likelihood estimation problem in an efficient and iterative manner. We will discuss more details in the following section, and here we will state the algorithm and give some examples.

One might wonder what would be the point of using an iterative algorithm for PCA while we already have a closed-form solution (see Proposition 4.4.2). One of the main point here is the computational efficiency. Namely, suppose that our data matrix  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$  has dimension  $d \times n$ , so each observation has dimension  $d$ . Say we are interested in learning  $r$  principal components, where  $r \ll d$  (e.g.,  $r = 25$  and  $d = 10^5$ ). The closed form solution in Proposition 4.4.2 involves computing the full  $d \times d$  covariance

matrix  $\mathbf{S} = n^{-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{X}})^T (\mathbf{x}_i - \bar{\mathbf{X}})$  and computing *all* of its eigenvector-eigenvalue pairs and keeping only the top  $r$  of them. However, computing the full  $d \times d$  covariance matrix takes  $O(nd^2)$  operations and computing its full eigendecomposition takes  $O(d^3)$  operations. But since we are only interested in extracting top  $r \ll d$  principal components, it would be much more efficient if we can avoid computing the full covariance matrix and still compute the top  $r$  principal components.

In Algorithm 4, we state the EM algorithm for PCA due to [Row98]. The computational cost of this algorithm is  $O(rdn)$ , which is much less compared to the computational cost  $O((n+d)d^2)$  of the full PCA solution.

---

**Algorithm 4** EM algorithm for PCA [Row98, TB99]

---

```

1: Input:  $\mathbf{X} \in \mathbb{R}^{d \times n}$  (data matrix);  $r \geq 1$  (rank parameter);  $\mathbf{W}_0 \in \mathbb{R}^{d \times r}$  (initial estimate);  $N$  (number of iterations)
2:   for  $n = 1, \dots, N$  do:
3:      $\mathbf{H}_n \leftarrow (\mathbf{W}_{n-1}^T \mathbf{W}_{n-1})^{-1} \mathbf{W}_{n-1}^T \mathbf{X} \in \mathbb{R}^{r \times n}$       (E-step)
4:      $\mathbf{W}_n \leftarrow \mathbf{X} \mathbf{H}_n^T (\mathbf{H}_n \mathbf{H}_n^T)^{-1} \in \mathbb{R}^{d \times r}$            (M-step)
5:   end for
6: output:  $\mathbf{W}_N$ 

```

---

Here are some important properties of Algorithm 4.

1. (*Convergence to stationary points*)  $\mathbf{W}_n$  converges to a stationary point of the PCA optimization problem (257) (and also the equivalent formulation (262)), but not guaranteed to converge to a global optimum.
2. (*Stability of global optimum*) If the initial estimate  $\mathbf{W}_0$  is close to a global optimum of the PCA problem  $\mathbf{W}^*$ , then  $\mathbf{W}_n$  converges to  $\mathbf{W}^*$ .
3. (*Instability of sub-optimal solutions*) If the initial estimate  $\mathbf{W}_0$  is close to a stationary point  $\mathbf{W}_*$  of the PCA problem and if  $\mathbf{W}_*$  is not a global optimum, then  $\mathbf{W}_n$  does not converge to  $\mathbf{W}_*$ .

Item 1 is due to the general convergence property of EM algorithm for maximum likelihood estimation problems. We will discuss more details for this aspect in later subsections. On the other hand, the local stability properties of stationary points of the PCA problem in items 2 and 3 are due to [TB99]. These are desired properties for iterative algorithms for nonconvex optimization problems, but they are not sufficient to guarantee convergence to a globally optimal solution (e.g., principal components for the PCA problem)<sup>13</sup>. For this reason, it is a general wisdom for EM algorithms to chose a proper initialization or try multiple runs for different initializations.

**6.2. Olivetti Face images.** Before we delve into theoretical details of EM algorithms for PCA, we first discuss some applications of it to the Olivetti face images. In Figure 60, we show the 24 components in  $\mathbf{W}_N$  learned by the EM-PCA in Algorithm 4 (100 trials each with random initialization  $\mathbf{W}_0$  and 100 iterations).

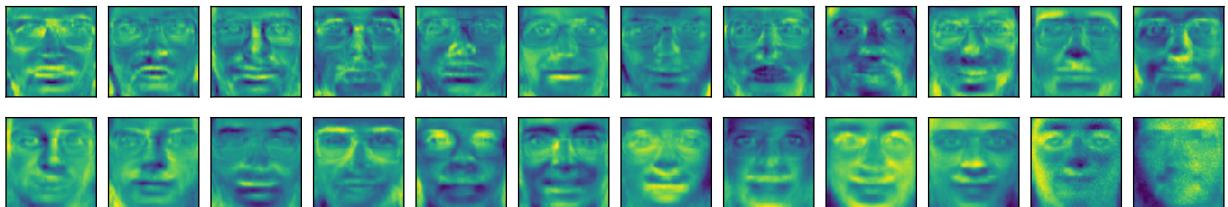


FIGURE 60. Examples of 24 dictionary images (principal components) learned from the Olivetti face dataset by the EM algorithm for PCA (Algorithm 4).

<sup>13</sup>In fact, the existence of a closed-form solution to PCA despite it being a non-convex problem is one of the reasons for the popularity of PCA.

The learned components are similar to the principal components in Figure 57, and if converged properly to a global optimum, these components should span a principal eigenspace of the covariance matrix  $\mathbf{S}$  of the input data  $\mathbf{X}$ <sup>14</sup>. However, in our experiment, it turns out that the components in Figure 57 are almost eigenvectors of  $\mathbf{S}$  but not quite. Indeed, in Figure 61, we see that the first column of  $\mathbf{W}_N$  learned by the EM-PCA is almost an eigenvector of  $\mathbf{S}$  but not as close as the exact principal component. In order to see this, we plot the ratio of the ( $64^2 = 4096$ ) coefficients in  $\mathbf{SW}[:, 0] / \mathbf{W}[:, 0]$ . If  $\mathbf{W}[:, 0]$  is an eigenvector of  $\mathbf{S}$ , then these ratios should all equal to some constant  $\lambda$ , which is the eigenvalue of  $\mathbf{W}[:, 0]$ . The standard deviation of these ratios are about 7.0567, while for the exact first principal component, the standard deviation is about  $7.7190 \times 10^{-5}$ .

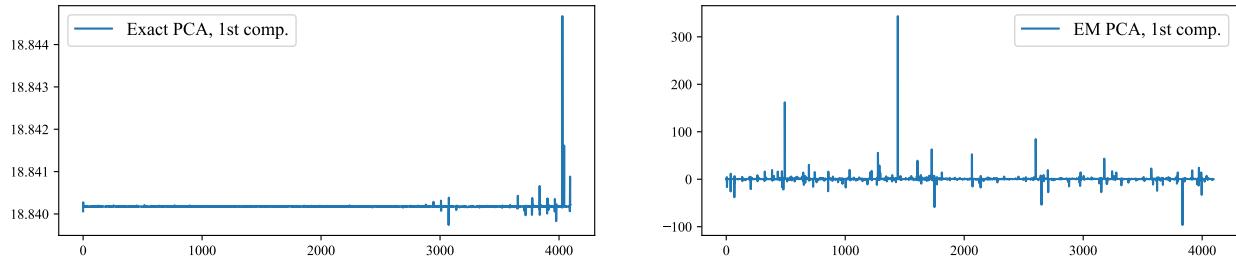


FIGURE 61. Plot of the ratio of coefficients  $\mathbf{SW}[:, 0] / \mathbf{W}[:, 0]$  where  $\mathbf{W}$  is either learned by the exact PCA or the EM PCA.  $\mathbf{W}[:, 0]$  is an eigenvector of  $\mathbf{S}$  if these ratios are constant.

**6.3. The EM algorithm for MLE.** The goal of the EM algorithm is to find maximum likelihood solutions for probabilistic models with latent variables. Suppose we have a probabilistic model for the pair  $(\mathbf{x}, \mathbf{h}) \in \mathbb{R}^p \times \mathbb{R}^r$  of random variables, given by a joint distribution  $(x, h) \mapsto p(x, h | \boldsymbol{\theta})$ , where  $\boldsymbol{\theta}$  denotes the vector of all parameters to be learned. We will refer to  $\mathbf{x}$  as the ‘observed variable’ and  $\mathbf{h}$  as the ‘latent variable’. Suppose we are given with observations  $x_1, \dots, x_n$ , which we will denote as a single data matrices  $X := [x_1, \dots, x_n] \in \mathbb{R}^{p \times n}$ . Each observed data  $x_i$  may be associated with latent variable  $h_i$ . We will denote the matrix of all latent variables as  $H := [h_1, \dots, h_n] \in \mathbb{R}^{r \times n}$ . When these are random, we will use the notations  $\mathbf{X} := [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{p \times n}$  and  $\mathbf{H} := [\mathbf{h}_1, \dots, \mathbf{h}_n] \in \mathbb{R}^{r \times n}$ .

The goal is to find a parameter  $\boldsymbol{\theta}$  that maximizes the likelihood of observing  $\mathbf{X}$  under the model, which can be written as (by using iterated expectation; see Exercise A.3.3)

$$\mathbb{P}_{\boldsymbol{\theta}}(\mathbf{X} = X) = \mathbb{E}_{\mathbf{H}} [\mathbb{P}(\mathbf{X} = X | \mathbf{H}, \boldsymbol{\theta})] = \sum_H p(X, H | \boldsymbol{\theta}). \quad (292)$$

As always in MLE problems, we assume that the columns of  $\mathbf{X}$  and  $\mathbf{H}$  are i.i.d., and in this case they follow the joint distribution  $p(x, h | \boldsymbol{\theta})$ . We may still denote the joint distribution of  $(\mathbf{X}, \mathbf{H})$  by  $p$ .

Hence we seek to find the maximum likelihood parameter

$$\hat{\boldsymbol{\theta}} := \arg \max_{\boldsymbol{\theta}} \left( \log (\mathbb{E}_{\mathbf{H}} [\mathbb{P}(\mathbf{X} = X | \mathbf{H}, \boldsymbol{\theta})]) = \log \left( \sum_H p(X, H | \boldsymbol{\theta}) \right) \right). \quad (293)$$

The key difficulty in solving the above optimization problem is that the expectation  $\mathbb{E}_{\mathbf{H}}$  over the latent variable  $\mathbf{H}$  (i.e., a sum over  $H$ ) is inside the log, which does not easily simplify and result in complicated analytical expression.

For the moment, suppose that an ‘oracle’ told us the unknown (random) latent variable  $\mathbf{H}$  is in fact a particular value  $H$ . Then we do not have to average over  $\mathbf{H}$  and simply need to solve

$$\arg \max_{\boldsymbol{\theta}} (\log (\mathbb{P}(\mathbf{X} = X | H, \boldsymbol{\theta})) = \log (p(X, H | \boldsymbol{\theta}))). \quad (294)$$

<sup>14</sup>Still in this case the columns  $\mathbf{W}_N$  need not be orthonormal, but one may use Gram-Schmidt to find an ONB for the column span of  $\mathbf{W}_N$  and corresponding eigenvalues can be found by multiplying each basis vector to  $\mathbf{S}$ .

However, in reality  $\mathbf{H}$  is random so we may prefer to maximize the expectation  $E_{\mathbf{H}}[\log(\mathbb{P}(\mathbf{X} = X | \mathbf{H}, \boldsymbol{\theta}))]$  instead of solving the above one-point problem. Now note that the distribution of  $\mathbf{H}$  conditional on  $X$  is given by

$$\mathbf{H} \sim p(H | X, \theta). \quad (295)$$

Even though  $X$  is known at this point, the parameter  $\theta$  is not. In order to proceed, we pretend  $\theta$  equals our current best estimate  $\theta^{\text{old}}$ . Then we know the distribution of  $H$  given  $X$  and  $\theta = \theta^{\text{old}}$ , so a better approximation of the original problem (293) than (294) would be the following:

$$\boldsymbol{\theta}^{\text{new}} = \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{H}} \left[ \log (\mathbb{P}(\mathbf{X} = X | \mathbf{H}, \boldsymbol{\theta})) \mid X, \boldsymbol{\theta}^{\text{old}} \right] \quad (296)$$

$$= \arg \max_{\boldsymbol{\theta}} \left( \sum_H \log(\mathbb{P}(\mathbf{X} = X | H, \boldsymbol{\theta})) p(H | X, \boldsymbol{\theta}^{\text{old}}) \right). \quad (297)$$

We can view (296) as an iterative algorithm that allows us to update the current best parameter  $\theta^{\text{old}}$  to a new one  $\theta^{\text{new}}$ . We call this the *expectation-maximization* (EM) algorithm since we first take the expectation  $E_H$  over the latent variable using the observed data  $X$  and the current best parameter  $\theta^{\text{old}}$ , and then maximize the resulting function in  $\theta$  to get the next estimate  $\theta^{\text{new}}$ .

**Algorithm 5** EM algorithm for MLE with latent variable

- ```

1: Input:  $\mathbf{X} \in \mathbb{R}^{d \times n}$  (data matrix);  $r \geq 1$  (rank parameter);  $\boldsymbol{\theta}_0 \in \mathbb{R}^{d \times r}$  (initial estimate);  $N$  (number of iterations)
2: for  $n = 1, \dots, N$  do
3:    $Q(\boldsymbol{\theta} | \boldsymbol{\theta}_{n-1}) \leftarrow \mathbb{E}_{\mathbf{H}} \left[ \log(\mathbb{P}(\mathbf{X} = X | \mathbf{H}, \boldsymbol{\theta})) \mid X, \boldsymbol{\theta}_{n-1} \right]$  (E-step)
4:    $\boldsymbol{\theta}_n \leftarrow \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta} | \boldsymbol{\theta}_{n-1})$  (M-step)
5: end for
6: output:  $\boldsymbol{\theta}_N$ 

```

**Proposition 4.6.1** (Monotonicity of EM algorithm). *Let  $\theta_n$  given by the EM algorithm (Algorithm 5). Then for all  $k \geq 0$ ,*

$$p(X|\boldsymbol{\theta}_k) \leq p(X|\boldsymbol{\theta}_{k+1}). \quad (298)$$

In other words, EM algorithm improves (increases) the likelihood  $p(X|\theta_k)$  monotonically.

PROOF. We will show the following inequality holds for all  $k \geq 1$ :

$$\log p(X|\boldsymbol{\theta}) - p(X|\boldsymbol{\theta}_k) \geq Q(\boldsymbol{\theta}|\boldsymbol{\theta}_k) - Q(\boldsymbol{\theta}_k|\boldsymbol{\theta}_k), \quad (299)$$

where  $Q(\boldsymbol{\theta} | \boldsymbol{\theta}_k) := \mathbb{E}_{\mathbf{H}} \left[ \log(p(X | \mathbf{H}, \boldsymbol{\theta})) \mid X, \boldsymbol{\theta}_k \right]$ . To see this, first note that  $p(X, H | \boldsymbol{\theta}) = p(X | H, \boldsymbol{\theta}) \cdot p(H | \boldsymbol{\theta})$ . So we can write

$$\log p(X|\theta) = \log p(X, H|\theta) - \log p(H|X, \theta) \quad (300)$$

Since the left hand side does not depend on  $H$ ,

$$\log p(X|\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{H}} \left[ \log p(X|\boldsymbol{\theta}) \mid X, \boldsymbol{\theta}_k \right]. \quad (301)$$

$$= \underbrace{\mathbb{E}_{\mathbf{H}} \left[ \log p(X, \mathbf{H} | \boldsymbol{\theta}) \middle| X, \boldsymbol{\theta}_k \right]}_{Q(\boldsymbol{\theta} | \boldsymbol{\theta}_k)} + \underbrace{\mathbb{E}_{\mathbf{H}} \left[ -\log p(\mathbf{H} | X, \boldsymbol{\theta}) \middle| X, \boldsymbol{\theta}_k \right]}_{H(\boldsymbol{\theta} | \boldsymbol{\theta}_k)}. \quad (302)$$

The above equation holds for all  $\theta$ , including  $\theta = \theta_k$ . So subtracting the two resulting equations,

$$\log p(X|\boldsymbol{\theta}) - \log p(X|\boldsymbol{\theta}_k) = Q(\boldsymbol{\theta}|\boldsymbol{\theta}_k) - Q(\boldsymbol{\theta}_k|\boldsymbol{\theta}_k) + [H(\boldsymbol{\theta}|\boldsymbol{\theta}_k) - H(\boldsymbol{\theta}_k|\boldsymbol{\theta}_k)]. \quad (303)$$

It then remains to show the last term in the bracket is nonnegative. Indeed,

$$H(\boldsymbol{\theta} | \boldsymbol{\theta}_k) - H(\boldsymbol{\theta}_k | \boldsymbol{\theta}_k) = \mathbb{E}_{\mathbf{H}} \left[ -\log \frac{p(\mathbf{H} | X, \boldsymbol{\theta})}{p(\mathbf{H} | X, \boldsymbol{\theta}_k)} \mid X, \boldsymbol{\theta}_k \right] \quad (304)$$

$$= D_{KL} \left( p(\cdot | X, \boldsymbol{\theta}) \parallel p(\cdot | X, \boldsymbol{\theta}_k) \right) \geq 0, \quad (305)$$

where the last equality uses the definition of KL divergence and the inequality is due to Gibb's inequality (see Exercise 4.6.3). This shows the assertion.  $\square$

**Exercise 4.6.2** (log-sum inequality). Let  $a_1, \dots, a_n \geq 0$  and  $b_1, \dots, b_n \geq 0$ . Show the following inequality holds:

$$\sum_{i=1}^n a_i \log \frac{a_i}{b_i} \geq \left( \sum_{i=1}^n a_i \right) \log \frac{\sum_{i=1}^n a_i}{\sum_{i=1}^n b_i}. \quad (306)$$

Furthermore, show that the equality above holds if and only if  $a_i/b_i$  does not depend on  $i = 1, \dots, n$ .

**Exercise 4.6.3** (Gibb's inequality). Let  $\mathbf{p}, \mathbf{q}$  denote probability distributions on the same probability space. Define the relative entropy (a.k.a. KullbackLeibler divergence) from  $\mathbf{q}$  to  $\mathbf{p}$  as

$$D_{KL}(\mathbf{p} \| \mathbf{q}) = \mathbb{E}_{X \sim \mathbf{p}} [-\log(\mathbf{q}(X)/\mathbf{p}(X))]. \quad (307)$$

(i) Show that

$$\mathbb{E}_{X \sim \mathbf{p}} [-\log(\mathbf{q}(X)/\mathbf{p}(X))] \geq -\log \mathbb{E}_{X \sim \mathbf{p}} [(\mathbf{q}(X)/\mathbf{p}(X))] \quad (308)$$

$$= -\log \int q(X) dX = -\log 1 = 0. \quad (309)$$

(ii) Deduce Gibb's inequality from (i):  $D_{KL}(\mathbf{p} \| \mathbf{q}) \geq 0$ , where the equality holds iff  $\mathbf{p} = \mathbf{q}$ .

(iii) Deduce the same result as in (ii) by using the log-sum inequality (see Exercise 4.6.2).

**6.4. Probabilistic PCA.** The *probabilistic PCA* [Row98, TB99] is the following *latent variable* model that explains a random observable  $\mathbf{x} \in \mathbb{R}^p$  as a linear transform of some lower dimensional latent variable  $\mathbf{h} \in \mathbb{R}^r$ :

$$\mathbf{x} = \mathbf{W}\mathbf{h} + \boldsymbol{\mu} + \boldsymbol{\epsilon}, \quad \mathbf{h} \sim N(\mathbf{0}, \mathbf{I}), \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}), \quad (310)$$

where  $\mathbf{h}$  and  $\boldsymbol{\epsilon}$  are independent. More explicitly,  $\mathbf{x}$  is modeled as a Gaussian RV:

$$\mathbf{x} \sim N(\boldsymbol{\mu}, \mathbf{WW}^T + \sigma^2 \mathbf{I}). \quad (311)$$

Indeed, it is easy to see that  $\mathbf{x}$  is a Gaussian RV, so we only need to compute its mean and covariance matrix as follow:

$$\mathbb{E}[\mathbf{x}] = \mathbb{E}[\mathbf{W}\mathbf{h} + \boldsymbol{\mu} + \boldsymbol{\epsilon}] = \mathbb{E}[\mathbf{W}\mathbf{h}] + \boldsymbol{\mu} = \mathbf{W}\mathbb{E}[\mathbf{h}] + \boldsymbol{\mu} = \boldsymbol{\mu}, \quad (312)$$

$$\text{Cov}(\mathbf{x}) = \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T] \quad (313)$$

$$= \mathbb{E}[(\mathbf{W}\mathbf{h} + \boldsymbol{\epsilon})(\mathbf{W}\mathbf{h} + \boldsymbol{\epsilon})^T] \quad (314)$$

$$= \mathbf{W}\mathbb{E}[\mathbf{h}\mathbf{h}^T]\mathbf{W}^T + \mathbb{E}[\boldsymbol{\epsilon}]\mathbb{E}[\mathbf{h}^T]\mathbf{W}^T + \mathbf{W}\mathbb{E}[\mathbf{h}]\mathbb{E}[\boldsymbol{\epsilon}^T] = \mathbf{WW}^T, \quad (315)$$

where for the covariance computation, we have used the fact that  $\mathbb{E}[\mathbf{h}\mathbf{h}^T] = \text{Cov}(\mathbf{h}) = \mathbf{I}$ , the independence between  $\mathbf{h}$  and  $\boldsymbol{\epsilon}$ , and  $\mathbb{E}[\boldsymbol{\epsilon}] = \mathbf{0}$ .

Denote  $\tilde{X} := X - [\boldsymbol{\mu}, \dots, \boldsymbol{\mu}]$  and  $\mathbf{C} := \mathbf{WW}^T + \sigma^2 \mathbf{I}$ . The log likelihood of observing  $X = [\vec{x}_1, \dots, \vec{x}_n]$  from (311) is

$$\ell(X | \boldsymbol{\mu}, \mathbf{W}, \sigma) = \sum_{i=1}^n \log(2\pi)^{-d/2} |\mathbf{C}|^{-1/2} - \frac{1}{2} (\vec{x}_i - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\vec{x}_i - \boldsymbol{\mu}) \quad (316)$$

$$= -\frac{n}{2} (d \log(2\pi) + \log |\mathbf{C}| + \text{tr}(\tilde{X} \mathbf{C}^{-1} \tilde{X})). \quad (317)$$

From this, by differentiating  $\mu$  and setting it equal zero, we find the maximum likelihood estimate of  $\mu$  is

$$\hat{\mu} = \bar{X} = \frac{1}{n} \sum_{i=1}^n \vec{x}_i. \quad (318)$$

On the other hand, the maximum likelihood estimates of  $\mathbf{W}$  and  $\sigma$  may be found by using EM algorithm. In order to do this, we first need the conditional distribution of the hidden variable  $\mathbf{h}$  given the data  $\mathbf{x}$  as in the following exercise.

**Exercise 4.6.4.** Let  $(\mathbf{x}, \mathbf{h})$  denote the pair of observable and latent variables in the probabilistic PCA in (310). Show that the distribution of  $\mathbf{h}$  given  $\mathbf{x}$  is the following Gaussian distribution:

$$\mathbf{h} | \mathbf{x} \sim N(\mathbf{M}^{-1} \mathbf{W}^T (\mathbf{x} - \mu), \sigma^{-2} \mathbf{M}^{-1}), \quad (319)$$

where  $\mathbf{M} := \mathbf{W}^T \mathbf{W} + \sigma^2 \mathbf{I} \in \mathbb{R}^{r \times r}$ . In particular, the expected value of  $\mathbf{h}$  given  $\mathbf{x}$  is

$$\mathbb{E}[\mathbf{h} | \mathbf{x}] = (\mathbf{W}^T \mathbf{W} + \sigma^2 \mathbf{I})^{-1} \mathbf{W}^T (\mathbf{x} - \mu). \quad (320)$$

(Hint: Use Bayes' Theorem and properties of multivariate Gaussian distribution, see [Bis06, Sec. 2.3.3].)

Now we are ready to derive the EM algorithm for probabilistic PCA.

**Proposition 4.6.5** (EM algorithm for PPCA). *Fix  $X \in \mathbb{R}^{d \times n}$  and denote  $\tilde{X} = X - [\bar{X}, \dots, \bar{X}]$  where  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X[:, i]$ . The EM algorithm for PPCA (Algorithm 5) equals the following iterative algorithm:*

$$\begin{cases} \mathbb{E}[\mathbf{H}_k] & \leftarrow (\mathbf{W}_{k-1}^T \mathbf{W}_{k-1} + \sigma_{k-1}^2 \mathbf{I}) \mathbf{W}_{k-1}^T \tilde{X}, \\ \mathbb{E}[\mathbf{H}_k \mathbf{H}_k^T] & \leftarrow \sigma_{k-1}^2 (\mathbf{W}_{k-1}^T \mathbf{W}_{k-1} + \sigma_{k-1}^2 \mathbf{I}) + \mathbb{E}[\mathbf{H}_k] \mathbb{E}[\mathbf{H}_k]^T \\ \mathbf{W}_k & \leftarrow \tilde{X} \mathbb{E}[\mathbf{H}_k^T] (\mathbb{E}[\mathbf{H}_k \mathbf{H}_k^T])^{-1}, \\ \sigma_k^2 & \leftarrow n^{-1} \text{tr}((\tilde{X} - \mathbf{W}_k \mathbb{E}[\mathbf{H}_k])^T (\tilde{X} - \mathbf{W}_k \mathbb{E}[\mathbf{H}_k])). \end{cases} \quad (321)$$

PROOF. (Ref: [TB99]) For the  $E$ -step, first we compute the conditional likelihood of observing the data  $X = [\vec{x}_1, \dots, \vec{x}_n]$  given the latent variable  $H = [\vec{h}_1, \dots, \vec{h}_n]$ . Denote  $\boldsymbol{\theta} = [\mathbf{W}, \sigma]$  and suppose we are given  $\boldsymbol{\theta}_{k-1} = [\mathbf{W}_{k-1}, \sigma_{k-1}]$ . Then

$$\ell(X | H, \boldsymbol{\theta}) := \sum_{i=1}^n \log(p(\vec{x}_i | \vec{h}_i, \boldsymbol{\theta})) \quad (322)$$

$$= n \log(\sigma^{-1} (2\pi)^{-d/2}) - \frac{1}{2\sigma^2} \sum_{i=1}^n \|\vec{x}_i - \mathbf{W}\vec{h}_i - \mu\|_F^2 \quad (323)$$

$$= n \log(\sigma^{-1} (2\pi)^{-d/2}) - \frac{1}{2\sigma^2} \text{tr}((\tilde{X} - \mathbf{W}H)^T (\tilde{X} - \mathbf{W}H)) \quad (324)$$

$$= n \log(\sigma^{-1} (2\pi)^{-d/2}) - \frac{1}{2\sigma^2} [\text{tr}(\tilde{X} \tilde{X}^T) - 2\text{tr}(X H^T \mathbf{W}^T) + \text{tr}(\mathbf{W} H H^T \mathbf{W}^T)]. \quad (325)$$

The  $E$ -step in Algorithm 5 amounts computing the following expectation:

$$Q(\boldsymbol{\theta} | \boldsymbol{\theta}_{k-1}) = \mathbb{E}_{\mathbf{H}} [\ell(X | \mathbf{H}, \boldsymbol{\theta}) | X, \boldsymbol{\theta}_{k-1}] \quad (326)$$

$$= n \log(\sigma^{-1} (2\pi)^{-d/2}) - \frac{1}{2\sigma^2} [\text{tr}(\tilde{X} \tilde{X}^T) - 2\text{tr}(\tilde{X} \mathbb{E}[\mathbf{H}]^T \mathbf{W}^T) + \text{tr}(\mathbf{W} \mathbb{E}[\mathbf{H} \mathbf{H}^T] \mathbf{W}^T)], \quad (327)$$

where, using Exercise (4.6.4),

$$\mathbb{E}[\mathbf{H}] := \mathbb{E}_{\mathbf{H}} [\mathbf{H} | X, \boldsymbol{\theta}_{k-1}] = \mathbf{M}^{-1} \mathbf{W}_n^T \tilde{X} = (\mathbf{W}_{k-1}^T \mathbf{W}_{k-1} + \sigma_{k-1}^2 \mathbf{I}) \mathbf{W}_{k-1}^T \tilde{X} \quad (328)$$

$$\mathbb{E}[\mathbf{H} \mathbf{H}^T] := \mathbb{E}_{\mathbf{H}} [\mathbf{H} \mathbf{H}^T | X, \boldsymbol{\theta}_{k-1}] = \text{Cov}(\mathbf{H}) + \mathbb{E}[\mathbf{H}] \mathbb{E}[\mathbf{H}]^T = \sigma_{k-1}^2 (\mathbf{W}_{k-1}^T \mathbf{W}_{k-1} + \sigma_{k-1}^2 \mathbf{I}) + \mathbb{E}[\mathbf{H}] \mathbb{E}[\mathbf{H}]^T. \quad (329)$$

Now differentiating  $Q(\boldsymbol{\theta} | \boldsymbol{\theta}_{k-1})$  above with respect to  $\mathbf{W}$  and  $\sigma$  gives

$$\nabla_{\mathbf{W}} Q(\boldsymbol{\theta} | \boldsymbol{\theta}_{k-1}) = -\sigma^{-2} [\mathbf{W}\mathbb{E}[\mathbf{H}\mathbf{H}^T] - \tilde{X}\mathbb{E}[\mathbf{H}]^T], \quad (330)$$

$$\frac{\partial}{\partial \sigma} Q(\boldsymbol{\theta} | \boldsymbol{\theta}_{k-1}) = -\frac{n}{\sigma} + \frac{1}{\sigma^3} \text{tr}((\tilde{X} - \mathbf{W}\mathbb{E}[\mathbf{H}])^T(\tilde{X} - \mathbf{W}\mathbb{E}[\mathbf{H}])). \quad (331)$$

Setting these equal to zero, we get

$$\mathbf{W}_n = \tilde{X}\mathbb{E}[\mathbf{H}]^T(\mathbb{E}[\mathbf{H}\mathbf{H}^T])^{-1}, \quad (332)$$

$$\sigma_n^2 = \frac{1}{n} \text{tr}((\tilde{X} - \mathbf{W}\mathbb{E}[\mathbf{H}])^T(\tilde{X} - \mathbf{W}\mathbb{E}[\mathbf{H}])). \quad (333)$$

Denoting  $\mathbb{E}[\mathbf{H}_n] = \mathbb{E}[\mathbf{H}]$  and  $\mathbb{E}[\mathbf{H}_n\mathbf{H}_n^T] = \mathbb{E}[\mathbf{H}\mathbf{H}^T]$ , we obtain the claimed iterative algorithm.  $\square$

**Remark 4.6.6.** Consider the special case of ‘zero-noise limit’  $\sigma \searrow 0$ . Then taking  $\sigma_k \searrow 0$  for each  $k \geq 1$ , the EM algorithm for PPCA in Proposition 4.6.5 reduces to

$$\begin{cases} \mathbf{H}_k & \leftarrow (\mathbf{W}_{n-1}^T \mathbf{W}_{n-1}) \mathbf{W}_{n-1}^T \tilde{X}, \\ \mathbf{W}_k & \leftarrow \tilde{X} \mathbf{H}_k^T (\mathbf{H}_k \mathbf{H}_k^T)^{-1}, \end{cases} \quad (334)$$

which is exactly Algorithm 4 we gave earlier.

## Bibliography

- [Bis06] Christopher M Bishop, *Pattern recognition and machine learning*, Springer, 2006.
- [HK70] Arthur E Hoerl and Robert W Kennard, *Ridge regression: Biased estimation for nonorthogonal problems*, Technometrics **12** (1970), no. 1, 55–67.
- [Hoy04] Patrik O Hoyer, *Non-negative matrix factorization with sparseness constraints.*, Journal of machine learning research **5** (2004), no. 9.
- [LKV21] Hanbaek Lyu, Yacoub H Kureh, Joshua Vendrow, and Mason A Porter, *Learning low-rank latent mesoscale structures in networks*, arXiv preprint arXiv:2102.06984 (2021).
- [LNB20] Hanbaek Lyu, Deanna Needell, and Laura Balzano, *Online matrix factorization for markovian data and applications to network dictionary learning*, Journal of Machine Learning Research **21** (2020), no. 251, 1–49.
- [LS99] Daniel D. Lee and H. Sebastian Seung, *Learning the parts of objects by non-negative matrix factorization*, Nature **401** (1999), no. 6755, 788.
- [LXT<sup>+</sup>17] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein, *Visualizing the loss landscape of neural nets*, arXiv preprint arXiv:1712.09913 (2017).
- [Lyu20] Hanbaek Lyu, *Convergence of block coordinate descent with diminishing radius for nonconvex optimization*, arXiv preprint arXiv:2012.03503 (2020).
- [RHW86] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams, *Learning representations by back-propagating errors*, nature **323** (1986), no. 6088, 533–536.
- [Row98] Sam Roweis, *EM algorithms for pca and spca*, Advances in neural information processing systems (1998), 626–632.
- [SH94] Ferdinando S Samaria and Andy C Harter, *Parameterisation of a stochastic model for human face identification*, Proceedings of 1994 IEEE workshop on applications of computer vision, IEEE, 1994, pp. 138–142.
- [SSY18] Tao Sun, Yuejiao Sun, and Wotao Yin, *On markov chain gradient descent*, arXiv preprint arXiv:1809.04216 (2018).
- [TB99] Michael E Tipping and Christopher M Bishop, *Probabilistic principal component analysis*, Journal of the Royal Statistical Society: Series B (Statistical Methodology) **61** (1999), no. 3, 611–622.
- [Tib96] Robert Tibshirani, *Regression shrinkage and selection via the lasso*, Journal of the Royal Statistical Society: Series B (Methodological) **58** (1996), no. 1, 267–288.
- [ZL12] Songfeng Zheng and Weixiang Liu, *Functional gradient ascent for probit regression*, Pattern recognition **45** (2012), no. 12, 4428–4437.

## APPENDIX A

# Review of basic probability theory

## 1. Probability measure and probability space

Many things in life are uncertain. Can we ‘measure’ and compare such uncertainty so that it helps us to make a more informed decision? Probability theory provides a systematic way of doing so.

We begin with idealizing our situation. Let  $\Omega$  be a finite set, called *sample space*. This is the collection of all possible outcomes that we can observe (think of six sides of a die). We are going to perform some experiment on  $\Omega$ , and the outcome could be any subset  $E$  of  $\Omega$ , which we call an *event*. Let us denote the collection of all events  $E \subseteq \Omega$  by  $2^\Omega$ . A *probability measure* on  $\Omega$  is a function  $\mathbb{P} : 2^\Omega \rightarrow [0, 1]$  such that for each event  $E \subseteq \Omega$ , it assigns a number  $\mathbb{P}(E) \in [0, 1]$  and satisfies the following properties:

- (i)  $\mathbb{P}(\emptyset) = 0$  and  $\mathbb{P}(\Omega) = 1$ .
- (ii) If two events  $E_1, E_2 \subseteq \Omega$  are disjoint, then  $\mathbb{P}(E_1 \cup E_2) = \mathbb{P}(E_1) + \mathbb{P}(E_2)$ .

In words,  $\mathbb{P}(E)$  is our quantization of how likely it is that the event  $E$  occurs out of our experiment.

**Exercise A.1.1.** Let  $\mathbb{P}$  be a probability measure on sample space  $\Omega$ . Show the following.

- (i) Let  $E = \{x_1, x_2, \dots, x_k\} \subseteq \Omega$  be an event. Then  $\mathbb{P}(E) = \sum_{i=1}^k \mathbb{P}(\{x_i\})$ .
- (ii)  $\sum_{x \in \Omega} \mathbb{P}(\{x\}) = 1$ .

If  $\mathbb{P}$  is a probability measure on sample space  $\Omega$ , we call the pair  $(\Omega, \mathbb{P})$  a *probability space*. This is our idealized world where we can precisely measure the uncertainty of all possible events. Of course, there could be many (in fact, infinitely many) different probability measures on the same sample space.

**Exercise A.1.2** (coin flip). Let  $\Omega = \{H, T\}$  be a sample space. Fix a parameter  $p \in [0, 1]$ , and define a function  $\mathbb{P}_p : 2^\Omega \rightarrow [0, 1]$  by  $\mathbb{P}_p(\emptyset) = 0$ ,  $\mathbb{P}_p(\{H\}) = p$ ,  $\mathbb{P}_p(\{T\}) = 1 - p$ ,  $\mathbb{P}_p(\{H, T\}) = 1$ . Verify that  $\mathbb{P}_p$  is a probability measure on  $\Omega$  for each value of  $p$ .

A typical way of constructing a probability measure is to specify how likely it is to see each individual element in  $\Omega$ . Namely, let  $f : \Omega \rightarrow [0, 1]$  be a function that sums up to 1, i.e.,  $\sum_{x \in \Omega} f(x) = 1$ . Define a function  $\mathbb{P} : 2^\Omega \rightarrow [0, 1]$  by

$$\mathbb{P}(E) = \sum_{\omega \in E} f(\omega). \quad (335)$$

Then this is a probability measure on  $\Omega$ , and  $f$  is called a *probability distribution* on  $\Omega$ . For instance, the probability distribution on  $\{H, T\}$  we used to define  $\mathbb{P}_p$  in Exercise A.1.2 is  $f(H) = p$  and  $f(T) = 1 - p$ .

**Example A.1.3** (Uniform probability measure). Let  $\Omega = \{1, 2, \dots, m\}$  be a sample space and let  $\mathbb{P}$  be the *uniform probability measure* on  $\Omega$ , that is,

$$\mathbb{P}(\{x\}) = 1/m \quad \forall x \in \Omega. \quad (336)$$

Then for the event  $A = \{1, 2, 3\}$ , we have

$$\mathbb{P}(A) = \mathbb{P}(\{1\} \cup \{2\} \cup \{3\}) \quad (337)$$

$$= \mathbb{P}(\{1\}) + \mathbb{P}(\{2\}) + \mathbb{P}(\{3\}) \quad (338)$$

$$= \frac{1}{m} + \frac{1}{m} + \frac{1}{m} = \frac{3}{m} \quad (339)$$

Likewise, if  $A \subseteq \Omega$  is any event and if we let  $|A|$  denote the size (number of elements) of  $A$ , then

$$\mathbb{P}(A) = \frac{|A|}{m}. \quad (340)$$

For example, let  $\Omega = \{1, 2, 3, 4, 5, 6\}^2$  be the sample space of a roll of two fair dice. Let  $A$  be the event that the sum of two dice is 5. Then

$$A = \{(1, 4), (2, 3), (3, 2), (4, 1)\}, \quad (341)$$

so  $|A| = 4$ . Hence  $\mathbb{P}(A) = 4/36 = 1/9$ . ▲

## 2. Random variables

Given a finite probability space  $(\Omega, \mathbb{P})$ , a (discrete) *random variable* (RV) is any real-valued function  $X : \Omega \rightarrow \mathbb{R}$ . We can think of it as the outcome of some experiment on  $\Omega$  (e.g., height of a randomly selected friend). We often forget the original probability space and specify a RV by its *probability mass function* (PMF)  $f_X : \mathbb{R} \rightarrow [0, 1]$ ,

$$f_X(x) = \mathbb{P}(X = x) = \mathbb{P}(\{\omega \in \Omega \mid X(\omega) = x\}). \quad (342)$$

Namely,  $\mathbb{P}(X = x)$  is the likelihood that the RV  $X$  takes value  $x$ .

**Example A.2.1.** Say you win \$1 if a fair coin lands heads and lose \$1 if lands tails. We can set up our probability space  $(\Omega, \mathbb{P})$  by  $\Omega = \{H, T\}$  and  $\mathbb{P}$  = uniform probability measure on  $\Omega$ . The RV  $X : \Omega \rightarrow \mathbb{R}$  for this game is  $X(H) = 1$  and  $X(T) = -1$ . The PMF of  $X$  is given by  $f_X(1) = \mathbb{P}(X = 1) = \mathbb{P}(\{H\}) = 1/2$  and likewise  $f_X(-1) = 1/2$ .

**Exercise A.2.2.** Let  $(\Omega, \mathbb{P})$  be a probability space and  $X : \Omega \rightarrow \mathbb{R}$  be a RV. Let  $f_X$  be the PMF of  $X$ , that is,  $f_X(x) = \mathbb{P}(X = x)$  for all  $x$ . Show that  $f_X$  adds up to 1, that is,

$$\sum_x f_X(x) = 1, \quad (343)$$

where the summation runs over all numerical values  $x$  that  $X$  can take.

There are two useful statistics of a RV to summarize its two most important properties: It's average and uncertainty. First, if one has to guess the value of a RV  $X$ , what would be the best choice? It is the *expectation* (or mean) of  $X$ , defined as below:

$$\mathbb{E}[X] = \sum_x x \mathbb{P}(X = x). \quad (344)$$

**Exercise A.2.3.** For any RV  $X$  and real numbers  $a, b \in \mathbb{R}$ , show that

$$\mathbb{E}[aX + b] = a\mathbb{E}[X] + b. \quad (345)$$

On the other hand, say you play two different games where in the first game, you win or lose \$1 depending on a fair coin flip, and in the second game, you win or lose \$10. In both games, your expected winning is 0. But the two games are different in how much the outcome fluctuates around the mean. This notion of the fluctuation is captured by the following quantity called *variance*:

$$\text{Var}(X) = \mathbb{E}[(X - \mathbb{E}(X))^2]. \quad (346)$$

Namely, it is the expected squared difference between  $X$  and its expectation  $\mathbb{E}(X)$ .

**Exercise A.2.4.** Let  $X$  be a RV. Show that  $\text{Var}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2$ .

Here are some of the simplest and yet most important RVs.

**Exercise A.2.5 (Bernoulli RV).** A RV  $X$  is a *Bernoulli* variable with (success) probability  $p \in [0, 1]$  if it takes value 1 with probability  $p$  and 0 with probability  $1 - p$ . In this case we write  $X \sim \text{Bernoulli}(p)$ . Show that  $\mathbb{E}(X) = p$  and  $\text{Var}(X) = p(1 - p)$ .

**Exercise A.2.6** (Indicator variables). Let  $(\Omega, \mathbb{P})$  be a probability space and let  $E \subseteq \Omega$  be an event. The *indicator variable* of the event  $E$ , which is denoted by  $\mathbf{1}_E$ , is the RV such that  $\mathbf{1}_E(\omega) = 1$  if  $\omega \in E$  and  $\mathbf{1}_E(\omega) = 0$  if  $\omega \in E^c$ . Show that  $\mathbf{1}_E$  is a Bernoulli variable with success probability  $p = \mathbb{P}(E)$ .

**Example A.2.7** (Uniform RV). Let  $\Omega = \{x_1, x_2, \dots, x_N\}$  be a finite sample space of real numbers. A RV  $X$  is a *uniform* variable on  $\Omega$ , denoted as  $X \sim \text{Uniform}(\Omega)$ , if it takes each of the element in  $\Omega$  with equal probability. That is,

$$\mathbb{P}(X = x_i) = \frac{1}{n} = \frac{1}{|\Omega|} \quad \text{for all } 1 \leq i \leq n. \quad (347)$$

Let  $X \sim \text{Uniform}(\Omega)$  as above. Then

$$\mathbb{E}[X] = \sum_{i=1}^n x_i \cdot \frac{1}{n} = \frac{1}{n} \sum_{i=1}^n x_i. \quad (348)$$

Also,

$$\text{Var}(X) = \sum_{i=1}^n (x_i - \mathbb{E}[X])^2 \cdot \frac{1}{n} = \frac{1}{n} \sum_{i=1}^n (x_i - \mathbb{E}[X])^2. \quad (349)$$

The following exercise ties the expectation and the variance of a RV into a problem of finding a point estimator that minimizes the mean squared error.

**Exercise A.2.8** (Variance as minimum MSE). Let  $X$  be a RV. Let  $\hat{x} \in \mathbb{R}$  be a number, which we consider as a ‘guess’ (or ‘estimator’ in Statistics) of  $X$ . Let  $\mathbb{E}[(X - \hat{x})^2]$  be the *mean squared error* (MSE) of this estimation.

(i) Show that

$$\mathbb{E}[(X - \hat{x})^2] = \mathbb{E}[X^2] - 2\hat{x}\mathbb{E}[X] + \hat{x}^2 \quad (350)$$

$$= (\hat{x} - \mathbb{E}[X])^2 + \mathbb{E}[X^2] - \mathbb{E}[X]^2 \quad (351)$$

$$= (\hat{x} - \mathbb{E}[X])^2 + \text{Var}(X). \quad (352)$$

(ii) Conclude that the MSE is minimized when  $\hat{x} = \mathbb{E}[X]$  and the global minimum is  $\text{Var}(X)$ . In this sense,  $\mathbb{E}[X]$  is the ‘best guess’ for  $X$  and  $\text{Var}(X)$  is the corresponding MSE.

To define a discrete RV, it was enough to specify its PMF. For a continuous RV, *probability distribution function* (PDF) plays an analogous role of PMF. We also need to replace summation  $\sum$  with an integral  $\int dx$ . Namely,  $X$  is a *continuous RV* if there is a function  $f_X : \mathbb{R} \rightarrow [0, \infty)$  such that for any interval  $[a, b]$ , the probability that  $X$  takes a value from an interval  $(a, b)$  is given by integrating  $f_X$  over the interval  $(a, b)$ :

$$\mathbb{P}(X \in (a, b)) = \int_a^b f_X(x) dx. \quad (353)$$

The *cumulative distribution function* (CDF) of a RV  $X$  (either discrete or continuous), denoted by  $F_X$ , is defined by

$$F_X(x) = \mathbb{P}(X \leq x). \quad (354)$$

By definition of PDF we get

$$F_X(x) = \int_{-\infty}^x f_X(t) dt. \quad (355)$$

Conversely, PDFs can be obtained by differentiating corresponding CDFs.

**Exercise A.2.9.** Let  $X$  be a continuous RV with PDF  $f_X$ . Let  $a$  be a continuity point of  $f_X$ , that is,  $f_X$  is continuous at  $a$ . Show that  $F_X(x)$  is differentiable at  $x = a$  and

$$\left. \frac{dF_X}{dx} \right|_{x=a} = f_X(a). \quad (356)$$

The expectation of a continuous RV  $X$  with pdf  $f_X$  is defined by

$$\mathbb{E}(X) = \int_{-\infty}^{\infty} x f_X(x) dx, \quad (357)$$

and its variance  $\text{Var}(X)$  is defined by the same formula (346).

### 3. Conditional and iterated expectations

Let  $X, Y$  be discrete RVs. Recall that the expectation  $\mathbb{E}(X)$  is the ‘best guess’ on the value of  $X$  when we do not have any prior knowledge on  $X$ . But suppose we have observed that some possibly related RV  $Y$  takes value  $y$ . What should be our best guess on  $X$ , leveraging this added information? This is called the *conditional expectation of  $X$  given  $Y = y$* , which is defined by

$$\mathbb{E}[X|Y = y] = \sum_x x \mathbb{P}(X = x|Y = y). \quad (358)$$

This best guess on  $X$  given  $Y = y$ , of course, depends on  $y$ . So it is a function in  $y$ . Now if we do not know what value  $Y$  might take, then we omit  $y$  and  $\mathbb{E}[X|Y]$  becomes a RV, which is called the *conditional expectation of  $X$  given  $Y$* .

**Proposition A.3.1** (Iterated expectation). *Let  $X, Y$  be RVs. Then  $\mathbb{E}[X] = \mathbb{E}[\mathbb{E}[X|Y]]$ .*

PROOF. We only show the assertion for the case when both  $X$  and  $Y$  are discrete, and the general case can be argued similarly.

We are going to write the iterated expectation  $\mathbb{E}[\mathbb{E}[X|Y]]$  as a double sum and swap the order of summation (Fubini’s theorem, as always).

$$\mathbb{E}[\mathbb{E}[X|Y]] = \sum_y \mathbb{E}[X|Y = y] \mathbb{P}(Y = y) \quad (359)$$

$$= \sum_y \left( \sum_x x \mathbb{P}(X = x|Y = y) \right) \mathbb{P}(Y = y) \quad (360)$$

$$= \sum_y \sum_x x \mathbb{P}(X = x|Y = y) \mathbb{P}(Y = y) \quad (361)$$

$$= \sum_y \sum_x x \mathbb{P}(X = x, Y = y) \quad (362)$$

$$= \sum_x \sum_y x \mathbb{P}(Y = y|X = x) \mathbb{P}(X = x) \quad (363)$$

$$= \sum_x x \left( \sum_y \mathbb{P}(Y = y|X = x) \right) \mathbb{P}(X = x) \quad (364)$$

$$= \sum_x x \mathbb{P}(X = x) = \mathbb{E}(X). \quad (365)$$

□

**Remark A.3.2.** Here is an intuitive reason why the iterated expectation works. Suppose you want to make the best guess  $\mathbb{E}(X)$ . Pretending you know  $Y$ , you can improve your guess to be  $\mathbb{E}(X|Y)$ . Then you admit that you didn’t know anything about  $Y$  and averaged over all values of  $Y$ . The result is  $\mathbb{E}[\mathbb{E}[X|Y]]$ , and this should be the same best guess on  $X$  when we don’t know anything about  $Y$ .

All our discussions above hold for continuous RVs as well: We simply replace the sum by integral and PMF by PDF. To summarize how we compute the iterated expectations when we condition on discrete and continuous RV:

$$\mathbb{E}[\mathbb{E}[X|Y]] = \begin{cases} \sum_y \mathbb{E}[X|Y = y] \mathbb{P}(Y = y) & \text{if } Y \text{ is discrete} \\ \int_{-\infty}^{\infty} \mathbb{E}[X|Y = y] f_Y(y) dy & \text{if } Y \text{ is continuous.} \end{cases} \quad (366)$$

**Exercise A.3.3** (Iterated expectation for probability). Let  $X, Y$  be RVs.

- (i) For any  $x \in \mathbb{R}$ , show that  $\mathbb{P}(X \leq x) = \mathbb{E}[\mathbf{1}(X \leq x)]$ .
- (ii) By using iterated expectation, show that

$$\mathbb{P}(X \leq x) = \mathbb{E}[\mathbb{P}(X \leq x | Y)], \quad (367)$$

where the expectation is taken over for all possible values of  $Y$ .

**Example A.3.4.** Let  $Y \sim \text{Uniform}([0, 1])$  and  $X \sim \text{Binomial}(n, Y)$ . Then  $X|Y = y \sim \text{Binomial}(n, y)$  so  $\mathbb{E}[X|Y = y] = ny$ . Hence

$$\mathbb{E}[X] = \int_0^1 \mathbb{E}[X | Y = y] f_Y(y) dy = \int_0^1 ny dy = n/2. \quad (368)$$

▲

**Example A.3.5.** Let  $X_1 \sim \text{Exp}(\lambda_1)$  and  $X_2 \sim \text{Exp}(\lambda_2)$  be independent exponential RVs. We will show that

$$\mathbb{P}(X_1 < X_2) = \frac{\lambda_1}{\lambda_1 + \lambda_2} \quad (369)$$

using the iterated expectation. Using iterated expectation for probability,

$$\mathbb{P}(X_1 < X_2) = \int_0^\infty \mathbb{P}(X_1 < X_2 | X_1 = x_1) \lambda_1 e^{-\lambda_1 x_1} dx_1 \quad (370)$$

$$= \int_0^\infty \mathbb{P}(X_2 > x_1) \lambda_1 e^{-\lambda_1 x_1} dx_1 \quad (371)$$

$$= \lambda_1 \int_0^\infty e^{-\lambda_2 x_1} e^{-\lambda_1 x_1} dx_1 \quad (372)$$

$$= \lambda_1 \int_0^\infty e^{-(\lambda_1 + \lambda_2)x_1} dx_1 = \frac{\lambda_1}{\lambda_1 + \lambda_2}. \quad (373)$$

▲

#### 4. Conditional variance

As we have defined conditional expectation, we may define the variance of a RV  $Y$  given that another RV  $X$  takes a particular value. Recall that the (unconditioned) variance of  $Y$  is defined by

$$\text{Var}(Y) = \mathbb{E}[(Y - \mathbb{E}(Y))^2]. \quad (374)$$

Note that there are two places where we take expectation. Given  $X$ , we should improve both expectations so the *conditional variance of  $Y$  given  $X$*  is defined by

$$\text{Var}(Y | X) = \mathbb{E}[(Y - \mathbb{E}[Y | X])^2 | X]. \quad (375)$$

**Proposition A.4.1.** Let  $X$  and  $Y$  be RVs. Then

$$\text{Var}(Y | X) = \mathbb{E}[Y^2 | X] - \mathbb{E}[Y | X]^2. \quad (376)$$

PROOF. Using linearity of conditional expectation and the fact that  $\mathbb{E}[Y | X]$  is not random given  $X$ ,

$$\text{Var}(Y | X) = \mathbb{E}[Y^2 - 2Y\mathbb{E}[Y | X] + \mathbb{E}[Y | X]^2 | X] \quad (377)$$

$$= \mathbb{E}[Y^2 | X] - \mathbb{E}[2Y\mathbb{E}[Y | X] | X] + \mathbb{E}[\mathbb{E}[Y | X]^2 | X] \quad (378)$$

$$= \mathbb{E}[Y^2 | X] - \mathbb{E}[Y | X]\mathbb{E}[2Y | X] + \mathbb{E}[Y | X]^2\mathbb{E}[1 | X] \quad (379)$$

$$= \mathbb{E}[Y^2 | X] - 2\mathbb{E}[Y | X]^2 + \mathbb{E}[Y | X]^2 \quad (380)$$

$$= \mathbb{E}[Y^2 | X] - \mathbb{E}[Y | X]^2. \quad (381)$$

□

The following exercise explains in what sense the conditional expectation  $\mathbb{E}[Y|X]$  is the best guess on  $Y$  given  $X$ , and that the minimum possible mean squared error is exactly the conditional variance  $\text{Var}(Y|X)$ .

**Exercise A.4.2.** Let  $X, Y$  be RVs. For any function  $g : \mathbb{R} \rightarrow \mathbb{R}$ , consider  $g(X)$  as an estimator of  $Y$ . Let  $\mathbb{E}_Y[(Y - g(X))^2 | X]$  be the *mean squared error*.

(i) Show that

$$\mathbb{E}_Y[(Y - g(X))^2 | X] = \mathbb{E}_Y[Y^2 | X] - 2g(X)\mathbb{E}_Y[Y | X] + g(X)^2 \quad (382)$$

$$= (g(X) - \mathbb{E}_Y(Y | X))^2 + \mathbb{E}_Y[Y^2 | X] - \mathbb{E}_Y[Y | X]^2 \quad (383)$$

$$= (g(X) - \mathbb{E}_Y(Y | X))^2 + \text{Var}(Y | X). \quad (384)$$

(ii) Conclude that the mean squared error is minimized when  $g(X) = \mathbb{E}_Y(Y | X)$  and the global minimum is  $\text{Var}(Y | X)$ .

Next, we study how we can decompose the variance of  $Y$  by conditioning on  $X$ .

**Proposition A.4.3** (Law of total variance). *Let  $X$  and  $Y$  be RVs. Then*

$$\text{Var}(Y) = \mathbb{E}_X[\text{Var}(Y | X)] + \text{Var}(\mathbb{E}[Y | X]). \quad (385)$$

PROOF. Using previous result, iterated expectation, and linearity of expectation, we have

$$\text{Var}(Y) = \mathbb{E}(Y^2) - (\mathbb{E}(Y))^2 \quad (386)$$

$$= \mathbb{E}_X(\mathbb{E}(Y^2 | X)) - (\mathbb{E}_X(\mathbb{E}(Y | X)))^2 \quad (387)$$

$$= \mathbb{E}_X(\text{Var}(Y | X) + (\mathbb{E}(Y | X))^2) - (\mathbb{E}_X(\mathbb{E}(Y | X)))^2 \quad (388)$$

$$= \mathbb{E}_X(\text{Var}(Y | X) + [\mathbb{E}_X(\mathbb{E}(Y | X))^2 - (\mathbb{E}_X(\mathbb{E}(Y | X)))^2]) \quad (389)$$

$$= \mathbb{E}_X(\text{Var}(Y | X)) + \text{Var}(\mathbb{E}(Y | X)). \quad (390)$$

□

Here is a handwavy explanation of why the above is true. Given  $Y$ , we should measure the fluctuation of  $Y | X$  from the conditional expectation  $\mathbb{E}[Y | X]$ , and this is measured as  $\text{Var}(Y | X)$ . Since we don't know  $Y$ , we average over all  $Y$ , giving  $\mathbb{E}[\text{Var}(Y | X)]$ . But the reference point  $\mathbb{E}[Y | X]$  itself varies with  $Y$ , so we should also measure its own fluctuation by  $\text{Var}(\mathbb{E}[Y | X])$ . These fluctuations add up nicely like Pythagorean theorem because  $\mathbb{E}[Y | X]$  is an optimal estimator so that these two fluctuations are 'orthogonal'.

**Example A.4.4.** Let  $X \sim \text{Uniform}([0, 1])$  and  $Y \sim \text{Binomial}(n, X)$ . Since  $Y | X = x \sim \text{Binomial}(n, x)$ , we have  $\mathbb{E}[Y | X = x] = nx$  and  $\text{Var}(Y | X = x) = nx(1 - x)$ . Also, since  $X \sim \text{Uniform}([0, 1])$ , we have

$$\text{Var}(\mathbb{E}[Y | X]) = \text{Var}(nX) = \frac{n^2}{12}. \quad (391)$$

So by iterated expectation, we get

$$\mathbb{E}[Y] = \mathbb{E}_X[\mathbb{E}[Y | X]] = \int_0^1 nx dx = \frac{n}{2}. \quad (392)$$

On the other hand, by law of total variance,

$$\text{Var}(Y) = \mathbb{E}[\text{Var}(Y | X)] + \text{Var}(\mathbb{E}[Y | X]) \quad (393)$$

$$= \int_0^1 nx(1 - x) dx + \text{Var}(nX) \quad (394)$$

$$= n \left[ \frac{x^2}{2} - \frac{x^3}{3} \right]_0^1 + \frac{n^2}{12} = \frac{n^2}{12} + \frac{n}{6}. \quad (395)$$

▲

## 5. Joint probability distributions

Let  $X$  and  $Y$  be discrete RVs defined on sample spaces  $\Omega_1$  and  $\Omega_2$ , which we think of the outcome of two random experiments. We can think of the pair  $(X, Y)$  of the outcomes as a single observable that describes the two experiments jointly. This object is called a *random vector*, being a vector of random variables. Namely,  $(X, Y)$  is a vector-valued function  $\Omega_1 \times \Omega_2 \rightarrow \mathbb{R}^2$  defined on the joint sample space  $\Omega_1 \times \Omega_2$ . As for the PDF of a random variable, which is the probability that a RV takes a particular value, we can also ask the probability that the random vector  $(X, Y)$  takes a particular vector  $(x, y)$ . This is called the *joint PMF* of  $X$  and  $Y$ :

$$f_{X,Y}(x, y) = \mathbb{P}((X, Y) = (x, y)) = \mathbb{P}(X = x, Y = y). \quad (396)$$

Given a joint PMF of  $X$  and  $Y$ , we can determine the PMFs of  $X$  and  $Y$  simply by ‘integrating out’ the other variable. Namely, note that for any  $x \in \mathbb{R}$ ,

$$\mathbb{P}(X = x) = \sum_{y \in Y[\Omega_2]} \mathbb{P}(X = x, Y = y) = \sum_{y \in Y[\Omega_2]} f_{X,Y}(x, y), \quad (397)$$

where  $Y[\Omega_2]$  denotes the set of all possible values of  $Y$ . Hence, we sum  $f_{X,Y}(x, y)$  over all  $y$  to get  $f_X(x)$ ; likewise, we sum  $f_{X,Y}(x, y)$  over all  $x$  to get  $f_Y(y)$ . We also call the PMFs of each  $X$  and  $Y$  the *marginal PMFs* of  $f_{X,Y}$ .

**Proposition A.5.1.** *Let  $X, Y$  be discrete RVs. Then they are independent if and only if there are functions  $f$  and  $g$  such that*

$$\mathbb{P}(X = x, Y = y) = f(x)g(y) \quad (398)$$

for all possible values  $(x, y)$  for  $(X, Y)$ .

PROOF. If  $X$  and  $Y$  are independent, then

$$\mathbb{P}(X = x, Y = y) = \mathbb{P}(X = x)\mathbb{P}(Y = y) = f_X(x)f_Y(y). \quad (399)$$

Hence the conclusion holds. Conversely, suppose (398) holds for some functions  $f$  and  $g$ . Then

$$\mathbb{P}(X = x) = \sum_y \mathbb{P}(X = x, Y = y) = \sum_y f(x)g(y) = f(x) \sum_y g(y). \quad (400)$$

Since sum of the probabilities  $\mathbb{P}(X = x)$  is 1, we have

$$\left( \sum_x f(x) \right) \left( \sum_y g(y) \right) = 1. \quad (401)$$

Hence if we define  $\bar{f}(x) = f(x) / \sum_x f(x)$  and  $\bar{g}(y) = \sum_y g(y)$ , then

$$\mathbb{P}(X = x) = f(x) \sum_y g(y) = \frac{f(x) \sum_y g(y)}{(\sum_x f(x)) (\sum_y g(y))} = \bar{f}(x) \sum_y \bar{g}(y) = \bar{f}(x). \quad (402)$$

Similarly, we have  $\mathbb{P}(Y = y) = \bar{g}(y)$ . This yields

$$\mathbb{P}(X = x, Y = y) = f(x)g(y) = \frac{f(x)g(y)}{(\sum_x f(x)) (\sum_y g(y))} = \bar{f}(x)\bar{g}(y) = \mathbb{P}(X = x)\mathbb{P}(Y = y). \quad (403)$$

Since  $(x, y)$  was arbitrary, this shows  $X$  and  $Y$  are independent.  $\square$

Let  $X$  and  $Y$  be continuous RVs defined on sample spaces  $\Omega_1$  and  $\Omega_2$ . We say  $X$  and  $Y$  are *jointly continuous* if the random vector  $(X, Y)$  has a function  $f_{X,Y} : \mathbb{R}^2 \rightarrow [0, \infty)$  such that for each subset  $\mathcal{A} \subseteq \mathbb{R}^2$ ,

$$\mathbb{P}((X, Y) \in \mathcal{A}) = \int \int_{(x,y) \in \mathcal{A}} f_{X,Y}(x, y) dx dy. \quad (404)$$

For instance, if  $\mathcal{A}$  is a rectangle of the form  $\mathcal{A} = [a, b] \times [c, d]$ , then by Fubini’s theorem, (404) reads

$$\mathbb{P}(a \leq X \leq b, c \leq Y \leq d) = \int_a^b \int_c^d f_{X,Y}(x, y) dx dy. \quad (405)$$

Such a function  $f_{X,Y}$  is called the *joint PDF* of  $X$  and  $Y$ .

As for the joint PMFs, we can determine the PDF of  $X$  and  $Y$  simply by ‘integrating out’ the other variable. We also call the PMFs of each  $X$  and  $Y$  the *marginal PMFs* of  $f_{X,Y}$ .

**Proposition A.5.2.** *Let  $X$  and  $Y$  be continuous RVs with joint PDF  $f_{X,Y}$ . Then for any  $x, y \in \mathbb{R}$ , we have*

$$f_X(x) = \int_{-\infty}^{\infty} f_{X,Y}(x, y) dy, \quad f_Y(y) = \int_{-\infty}^{\infty} f_{X,Y}(x, y) dx. \quad (406)$$

PROOF. Note that for any interval  $[a, b] \subseteq \mathbb{R}$ , by definition of joint PDF and Fubini’s theorem,

$$\mathbb{P}(a \leq X \leq b) = \int_a^b \int_{-\infty}^{\infty} f_{X,Y}(x, y) dx dy \quad (407)$$

$$= \int_a^b \left( \int_{-\infty}^{\infty} f_{X,Y}(x, y) dy \right) dx. \quad (408)$$

Then by the definition of PDF of a continuous RV, it follows that the PDF  $f_X$  of  $X$  is given the first equation in (406). A Similar argument applies to  $f_Y$  as well.  $\square$

**Example A.5.3.** Let  $X$  and  $Y$  be continuous RVs with joint PDF given by

$$f_{X,Y}(x, y) = \lambda^2 e^{-\lambda(x+y)} \mathbf{1}(x, y \geq 0). \quad (409)$$

Then for each  $x \in \mathbb{R}$ ,

$$f_X(x) = \int_{-\infty}^{\infty} \lambda^2 e^{-\lambda(x+y)} \mathbf{1}(x, y \geq 0) dy \quad (410)$$

$$= \lambda e^{-\lambda x} \mathbf{1}(x \geq 0) \int_{-\infty}^{\infty} \lambda e^{-\lambda y} \mathbf{1}(y \geq 0) dy \quad (411)$$

$$= \lambda e^{-x} \mathbf{1}(x \geq 0). \quad (412)$$

Note that for the last equality, we have recognized the function  $y \mapsto \lambda e^{-\lambda y} \mathbf{1}(y \geq 0)$  as the PDF of  $\text{Exp}(\lambda)$  RV. This shows

$$f_X(x) = \lambda e^{-\lambda x} \mathbf{1}(x \geq 0). \quad (413)$$

Hence  $X \sim \text{Exp}(\lambda)$ . Similarly, we obtain  $Y \sim \text{Exp}(\lambda)$ .  $\blacktriangle$

## 6. Definition and Examples of MLE

Maximum Likelihood Estimation (MLE) is a classical notion in statistics, which has a prominent use in the context of modern machine learning. A core problem in statistics is to infer an unknown RV  $X$  from its sample values  $x_1, \dots, x_n$ . In many cases, we can narrow down our inference problem by imposing a probabilistic (statistical) model for  $X$ . Namely, say we know  $X$  is a Bernoulli RV with unknown success probability  $p$ . We then only need to estimate the parameter  $p$  using our sample, not the entire distribution of  $X$ . MLE gives a systematic approach to this parameter estimation problem. Below we give the pipeline of MLE.

### Pipeline of MLE.

**Input:** An unknown RV  $X$  with distribution  $f_{X;\theta}$ , parameterized by  $\theta$  in a parameter space  $\Omega$ . Also have sample values  $x_1, x_2, \dots, x_n$ .

**Objective:** Obtain an estimation  $\hat{\theta}$  of  $\theta$  using the sample.

**Method:** Choose  $\hat{\theta} \in \Omega$  so that it maximizes the following *likelihood function*

$$L(x_1, \dots, x_n; \theta) = \prod_{i=1}^n f_{X,\theta}(x_i). \quad (414)$$

Here the RV  $\hat{\theta}(X_1, \dots, X_n)$  is called the *maximum likelihood estimator* of  $\theta$ , and its observed value  $\hat{\theta}(x_1, \dots, x_n)$  is called the *maximum likelihood estimate* of  $\theta$ .

What is the reasoning behind maximizing the likelihood function as above? The underlying assumption is that, *you are seeing the sample values  $x_1, \dots, x_n$  because it was the most likely to see them!* Namely, suppose we have designed the sampling procedure well-enough so that we get i.i.d. samples  $X_1, \dots, X_n$  each with distribution  $f_{X;\theta}$ . Then we have

$$L(x_1, \dots, x_n; \theta) = \mathbb{P}(X_1 = x_1, \dots, X_n = x_n; \theta). \quad (415)$$

Namely, the likelihood function on the LHS is the probability of observing a sequence of specific sample values  $x_1, \dots, x_n$  under the i.i.d. assumption and assuming the parameter value  $\theta$ . Hence, MLE chooses  $\hat{\theta}$  to be the value of the parameter in  $\Omega$  under which the probability of obtaining the current sample is maximized.

**Example A.6.1** (MLE for Bernoulli RVs). Suppose  $X \sim \text{Bernoulli}(p)$  for some unknown  $p \in [0, 1]$ . Suppose we have sample values  $x_1, \dots, x_n$  obtained from an i.i.d. sampling for  $X$ . Denote  $\bar{x} = n^{-1} \sum_{i=1}^n x_i$ . We will show that the MLE for  $p$  is  $\bar{X}$ , that is,

$$\hat{p} = \bar{X}. \quad (416)$$

Let  $X_1, \dots, X_n$  be i.i.d. with distribution  $\text{Bernoulli}(p)$ . Note that

$$\mathbb{P}(X_i = x_i; p) = \begin{cases} p & \text{if } x_i = 1 \\ 1-p & \text{if } x_i = 0 \end{cases} \quad (417)$$

$$= p^{x_i} (1-p)^{1-x_i}. \quad (418)$$

Hence the likelihood function is given by

$$L(x_1, \dots, x_n; p) = \prod_{i=1}^n \mathbb{P}(X_i = x_i; p) \quad (419)$$

$$= \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} \quad (420)$$

$$= p^{\sum_{i=1}^n x_i} (1-p)^{n - \sum_{i=1}^n x_i}. \quad (421)$$

In order to maximize the likelihood function, it suffices to maximize its logarithm <sup>1</sup>. The *log likelihood function* is given by

$$l(x_1, \dots, x_n; p) := \log L(x_1, \dots, x_n; p) = n\bar{x}\log p + (n - n\bar{x})\log(1-p). \quad (422)$$

To maximize the log likelihood function, we take partial derivative in  $p$ :

$$\frac{\partial l(x_1, \dots, x_n; p)}{\partial p} = \frac{n\bar{x}}{p} - \frac{n - n\bar{x}}{1-p}. \quad (423)$$

Setting this equal to zero, we obtain

$$\frac{\bar{x}}{p} = \frac{1 - \bar{x}}{1-p}. \quad (424)$$

Rearranging, we obtain  $p = \bar{x}$ . Hence we have (416) as desired. ▲

**Example A.6.2** (MLE for Exponential RVs). Suppose  $X \sim \text{Exp}(\lambda)$  for some unknown  $\lambda > 0$ . We have sample values  $x_1, \dots, x_n$  obtained from an i.i.d. sampling for  $X$ . Denote  $\bar{x} = n^{-1} \sum_{i=1}^n x_i$ . We will show that the MLE for  $\lambda$  is  $1/\bar{X}$ , that is,

$$1/\hat{\lambda} = \bar{X}. \quad (425)$$

---

<sup>1</sup>since  $\log$  is a strictly increasing function.

Let  $X_1, \dots, X_n$  be i.i.d. with distribution  $\text{Exp}(\lambda)$ , which have the following PDF

$$f_X(x; \lambda) = \lambda e^{-\lambda x} \mathbf{1}(x \geq 0). \quad (426)$$

Noting that  $x_1, \dots, x_n \geq 0$ , it follows that the likelihood function is given by

$$L(x_1, \dots, x_n; \lambda) = \prod_{i=1}^n f_X(x_i; \lambda) \quad (427)$$

$$= \lambda^n \prod_{i=1}^n e^{-\lambda x_i} = \lambda^n e^{-\lambda \sum_{i=1}^n x_i} = \lambda^n e^{-\lambda n \bar{x}}. \quad (428)$$

So the log likelihood function is given by

$$l(x_1, \dots, x_n; \lambda) = n \log \lambda - \lambda n \bar{x}. \quad (429)$$

To maximize the log likelihood function, we take partial derivative in  $\lambda$ :

$$\frac{\partial l(x_1, \dots, x_n; \lambda)}{\partial \lambda} = \frac{n}{\lambda} - n \bar{x}. \quad (430)$$

Setting this equal to zero, we obtain  $1/\lambda = \bar{x}$ . Hence  $1/\hat{\lambda} = \bar{X}$ , as desired.  $\blacktriangle$

**Example A.6.3** (MLE for Normal RVs). Suppose  $X \sim N(\mu, \sigma^2)$  for some unknown  $\mu \in \mathbb{R}$  and  $\sigma \geq 0$ . Let  $X_1, \dots, X_n$  be i.i.d. with distribution  $N(\mu, \sigma^2)$ . Denote the sample mean and variance of the empirical distribution by

$$\bar{X} = n^{-1} \sum_{i=1}^n X_i, \quad V = n^{-1} \sum_{i=1}^n (X_i - \bar{X})^2. \quad (431)$$

We will show that the MLE for  $\mu$  and  $\sigma^2$  are given by sample mean and variance of the empirical distribution:

$$\hat{\mu} = \bar{X}, \quad (\hat{\sigma}^2) = V. \quad (432)$$

Recall that  $N(\mu, \sigma^2)$  has the following PDF

$$f_X(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right). \quad (433)$$

It follows that the likelihood function is given by

$$L(x_1, \dots, x_n; \mu, \sigma^2) = \prod_{i=1}^n f_X(x_i; \mu, \sigma^2) \quad (434)$$

$$= (2\pi\sigma^2)^{-n/2} \prod_{i=1}^n \exp\left(-\frac{(x_i-\mu)^2}{2\sigma^2}\right) \quad (435)$$

$$= (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i-\mu)^2\right). \quad (436)$$

So the log likelihood function is given by

$$l(x_1, \dots, x_n; \mu, \sigma^2) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i-\mu)^2. \quad (437)$$

To maximize the log likelihood function, we take partial derivatives in  $\mu$  and  $\sigma^2$ :

$$\frac{\partial l(x_1, \dots, x_n; \mu, \sigma^2)}{\partial \mu} = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = \frac{n}{\sigma^2} (\bar{x} - \mu), \quad (438)$$

$$\frac{\partial l(x_1, \dots, x_n; \mu, \sigma^2)}{\partial \sigma^2} = -\frac{n}{2} \frac{1}{\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2. \quad (439)$$

Setting these equations to be zero, we obtain

$$\begin{cases} \bar{x} = \mu \\ -n + \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 = 0. \end{cases} \quad (440)$$

Using the first equation for the second, we obtain

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2. \quad (441)$$

This shows (432). ▲

## 7. Bayesian Inference

In this section, we discuss another widely used method of parameter estimation, called the *Bayesian estimation*. Just like in MLE, we would like to estimate an unknown parameter  $\theta$  in a distribution  $f_{X;\theta}$  from analyzing an i.i.d. sample  $X_1, \dots, X_n$  drawn from that distribution. One of the key aspects in Bayesian estimation is to quantify our prior knowledge on the unknown parameter  $\theta$  as yet another probabilistic model, which we call as *belief*. After observing a new *data*  $\{X_1 = x_1, \dots, X_n = x_n\}$ , our current best belief on  $\theta$  (*prior distribution*) will be updated to a new belief (posterior distribution) according to Bayes' Theorem. Since we are making some additional modeling assumptions on  $\theta$ , we will be able to draw more quantitative inference on  $\theta$  using Bayesian estimation.

One of the early bottlenecks in Bayesian estimation was the computational overhead in computing posterior distributions. Due to the recent explosion in our computational capacity, coupled with advanced sampling techniques such as MCMC, Bayesian methods have become one of the indispensable tools in modern statistics and machine learning.

We begin by recalling the following basic observation, from which Bayes' theorem can be easily derived:

$$\mathbb{P}(B|A)\mathbb{P}(A) = \mathbb{P}(A \cap B) = \mathbb{P}(A|B)\mathbb{P}(B). \quad (442)$$

In the case of the events involving continuous random variables taking a certain value, we interpret the above probabilities as probability densities. The following is a random variable version of Bayes' Theorem discussed in the previous subsection.

**Theorem A.7.1** (Bayes' Theorem). *Let  $X$  and  $\Theta$  be random variables. Then*

$$\overbrace{\mathbb{P}(\Theta = \theta | X = x)}^{\text{posterior}} = \frac{\mathbb{P}(X = x | \Theta = \theta)\mathbb{P}(\Theta = \theta)}{\mathbb{P}(X = x)} \propto \underbrace{\mathbb{P}(X = x | \Theta = \theta)}_{\text{likelihood}} \underbrace{\mathbb{P}(\Theta = \theta)}_{\text{prior}} \quad (443)$$

PROOF. Follows from (442). □

**Remark A.7.2.** Depending on whether  $X$  and  $\Theta$  are discrete or continuous, we interpret the probabilities in the above theorem as PDF or PMF, accordingly.

Bayesian inference is usually carried out in the following procedure.

### Bayesian inference:

- (i) To explain an observable  $\mathbf{x}$ , we choose a *probabilistic model*  $p(x|\theta)$ , which is a probability distribution on the possible values  $x$  of  $\mathbf{x}$ , depending on a parameter  $\theta$ .
- (ii) Choose a probability distribution  $\pi(\theta)$ , called the *prior distribution*, that expresses our beliefs about a parameter  $\theta$  to the best of our current knowledge.
- (iii) After observing data  $\mathcal{D} = \{x_1, \dots, x_n\}$ , we update our beliefs and compute the *posterior distribution*  $p(\theta|\mathcal{D})$  according to Bayes' Theorem.

When we generate data  $\mathcal{D} = \{x_1, \dots, x_n\}$ , we assume that each  $x_i$  is obtained by an independent sample  $X_i$  of the observable  $\mathbf{x}$  from the true distribution of  $\mathbf{x}$ . According to Bayes' Theorem, we can compute the posterior distribution as

$$\mathbb{P}(\Theta = \theta | X_1 = x_1, \dots, X_n = x_n) = \frac{\mathbb{P}(X_1 = x_1, \dots, X_n = x_n | \Theta = \theta) \mathbb{P}(\Theta = \theta)}{\mathbb{P}(X_1 = x_1, \dots, X_n = x_n)}, \quad (444)$$

or in a more compact form,

$$p(\theta | x_1, \dots, x_n) = \frac{p(x_1, \dots, x_n | \theta) \pi(\theta)}{p(x_1, \dots, x_n)}. \quad (445)$$

By a slight abuse of notation, we use lowercase  $p$  to denote either PMF or PDF depending on the context.

The conditional probability  $p(x_1, \dots, x_n | \theta)$  is called the *likelihood function*, which is the probability of obtaining independent data sample  $\mathcal{D} = \{x_1, \dots, x_n\}$  according to our probability model  $p(x | \theta)$  assuming model parameter  $\Theta = \theta$ . By the independence between samples, the likelihood function factors into the product of marginal likelihood of each sample:

$$p(x_1, \dots, x_n | \theta) = \mathbb{P}(X_1 = x_1, \dots, X_n = x_n | \Theta = \theta) \quad (446)$$

$$= \prod_{i=1}^n \mathbb{P}(X_i = x_i | \Theta = \theta) = \prod_{i=1}^n p(x_i | \theta). \quad (447)$$

On the other hand, the joint probability  $\mathbb{P}(X_1 = x_1, \dots, X_n = x_n)$  of obtaining the data sample  $\mathcal{D} = \{x_1, \dots, x_n\}$  from our probability model can be computed by conditioning on the values of model parameter  $\Theta$ :

$$\mathbb{P}(X_1 = x_1, \dots, X_n = x_n) = \mathbb{E} [\mathbb{P}(X_1 = x_1, \dots, X_n = x_n | \Theta) | \Theta] \quad (448)$$

$$= \int_{-\infty}^{\infty} p(x_1, \dots, x_n | \theta) \pi(\theta) d\theta. \quad (449)$$

What's so important about Bayes' theorem is its interpretation as a means of *inference*, which is one of the fundamental tools in modern machine learning.

**Example A.7.3.** Suppose Bob has a coin with an unknown probability of heads, which we denote by  $\Theta$ . This is called the *parameter*. Suppose Alice knows that Bob has one of the three kinds of coins  $A$ ,  $B$ , and  $C$ , with the probability of heads being  $p_A = 0.2$ ,  $p_B = 0.5$ , and  $p_C = 0.8$ , respectively. This piece of information is called the *model*. Since Alice has no information, she initially assumes that Bob has one of the three coins equally likely. Namely, she assumes the uniform distribution over the sample space  $\Omega = \{0.2, 0.5, 0.8\}$ . This knowledge is called *prior*.

Now Bob flips his coin 10 times and got 7 heads, and reports this information, which we call Data, to Alice. Now that Alice has more information, she needs to update her *prior* to *posterior*, which is the probability distribution on  $\Omega$  that best explains the Data. Namely, it is the conditional probability distribution  $\mathbb{P}(\Theta = \theta | \text{Data})$ .

First, using our prior, we compute  $\mathbb{P}(\text{Data})$ , the probability of seeing this particular data at hand. This can be done by partitioning:

$$\mathbb{P}(\text{Data}) = \mathbb{P}(\text{Data} | \Theta = 0.2) \mathbb{P}(\Theta = 0.2) + \mathbb{P}(\text{Data} | \Theta = 0.5) \mathbb{P}(\Theta = 0.5) \quad (450)$$

$$+ \mathbb{P}(\text{Data} | \Theta = 0.8) \mathbb{P}(\Theta = 0.8) \quad (451)$$

$$= \mathbb{P}(7/10 \text{ heads} | \Theta = 0.2) \frac{1}{3} + \mathbb{P}(7/10 \text{ heads} | \Theta = 0.5) \frac{1}{3} + \mathbb{P}(7/10 \text{ heads} | \Theta = 0.8) \frac{1}{3} \quad (452)$$

$$= \frac{1}{3} \left( \binom{10}{7} (0.2)^7 (0.8)^3 + \binom{10}{7} (0.5)^7 (0.5)^3 + \binom{10}{7} (0.8)^7 (0.2)^3 \right) \approx 0.1064, \quad (453)$$

where  $\binom{10}{7} = 120$  is the number of ways to choose 7 out of 10 objects.

Second, we reformulate the first equality of Theorem A.7.1 as

$$\mathbb{P}(\Theta | \text{Data}) = \frac{\mathbb{P}(\text{Data} | \Theta) \mathbb{P}(\Theta)}{\mathbb{P}(\text{Data})}. \quad (454)$$

Hence we can compute the posterior distribution by

$$\mathbb{P}(\Theta = 0.2 | \text{Data}) = \frac{\mathbb{P}(\text{Data} | \Theta = 0.2) \mathbb{P}(\Theta = 0.2)}{\mathbb{P}(\text{Data})} = \frac{\binom{10}{7} (0.2)^7 (0.8)^3 \frac{1}{3}}{0.1064} \approx 0.0025 \quad (455)$$

$$\mathbb{P}(\Theta = 0.5 | \text{Data}) = \frac{\mathbb{P}(\text{Data} | \Theta = 0.5) \mathbb{P}(\Theta = 0.5)}{\mathbb{P}(\text{Data})} = \frac{\binom{10}{7} (0.5)^7 (0.5)^3 \frac{1}{3}}{0.1064} \approx 0.3670 \quad (456)$$

$$\mathbb{P}(\Theta = 0.8 | \text{Data}) = \frac{\mathbb{P}(\text{Data} | \Theta = 0.8) \mathbb{P}(\Theta = 0.8)}{\mathbb{P}(\text{Data})} = \frac{\binom{10}{7} (0.8)^7 (0.2)^3 \frac{1}{3}}{0.1064} \approx 0.6305. \quad (457)$$

Note that according to the posterior distribution,  $\Theta = 0.8$  is the most likely value, which is natural given that we have 7 heads out of 10 flips. However, our knowledge is always incomplete so our posterior knowledge is still a probability distribution on the sample space.

What if Bob flips his coin another 10 times and reports only 3 heads to Alice? Then she will have to use her current prior  $\pi = [0.0025, 0.3670, 0.6305]$  (which was obtained as the posterior in the previous round) to compute yet another posterior using the new data. This will be likely to give a higher weight to  $\Theta = 0.2$ .  $\blacktriangle$

**Exercise A.7.4.** Suppose we have a prior distribution  $\pi = [0.0025, 0.3670, 0.6305]$  on the sample space  $\Omega = \{0.2, 0.5, 0.8\}$  for the inference problem of unknown parameter  $\Theta$ . Suppose we are given the data that ten independent flips of probability  $\Theta$  coin come up heads twice. Compute the posterior distribution using this data and Bayesian inference.

**Exercise A.7.5.** A test for pancreatic cancer is assumed to be correct 95% of the time: if a person has cancer, the test results in positive with a probability 0.95, and if the person does not have cancer, then the test results in negative with probability 0.95. From recent medical research, it is known that only 5% of the population has pancreatic cancer. Given that the person just tested positive, what is the probability of having cancer?

**Example A.7.6** (Bernoulli model and uniform prior). Bob has a coin with unknown probability  $\Theta$  of heads. Alice has no information whatsoever, so her prior distribution  $\pi$  for  $\Theta$  is the uniform distribution  $\text{Uniform}([0, 1])$ . Bob flips his coin independently  $n$  times, and let  $X_1, \dots, X_n$  be the outcome,  $X_i$ 's are i.i.d. Bernoulli( $\Theta$ ) variables. Let  $x_i$  be the observed value of  $X_i$ , and let  $s_n = x_1 + \dots + x_n$  be the number of heads in the  $n$  flips. Given data  $\mathcal{D} = \{x_1, \dots, x_n\}$ , Alice wants to compute her posterior distribution on  $\Theta$ .

The likelihood function is given by

$$p(x_1, \dots, x_n | \theta) = \prod_{i=1}^n \theta^{x_i} (1-\theta)^{1-x_i} \quad (458)$$

$$= \theta^{s_n} (1-\theta)^{n-s_n}, \quad (459)$$

where we denote  $s_n = x_1 + \dots + x_n$ . Since  $\pi \equiv 1$ , the posterior distribution is given by

$$p(\theta | x_1, \dots, x_n) = \frac{\theta^{s_n} (1-\theta)^{n-s_n}}{p(x_1, \dots, x_n)}, \quad (460)$$

Note that

$$p(x_1, \dots, x_n) = \int_0^1 \theta^{s_n} (1-\theta)^{n-s_n} d\theta = \frac{1}{\binom{n}{s_n} (n+1)}, \quad (461)$$

where the last equality follows from Exercise A.7.8. Hence

$$p(\theta | x_1, \dots, x_n) = \binom{n}{s_n} (n+1) \theta^{s_n} (1-\theta)^{n-s_n} \quad (462)$$

$$= \frac{(n+1)!}{s_n!(n-s_n)!} \theta^{(s_n+1)-1} (1-\theta)^{(n-s_n+1)-1}. \quad (463)$$

Hence, according to Exercise A.7.9, we can write

$$\Theta | x_1, \dots, x_n \sim \text{Beta}(s_n + 1, n - s_n + 1). \quad (464)$$

▲

**Exercise A.7.7.** Bob has a coin with unknown probability  $\Theta$  of heads. Alice has the following Beta prior (See Exercise A.7.9 for the definition of Beta distribution):

$$\pi = \text{Beta}(\alpha, \beta). \quad (465)$$

Suppose that Bob gives Alice the data  $\mathcal{D}_n = \{x_1, \dots, x_n\}$ , which is the outcome of  $n$  independent coin flips. Denote  $s_n = x_1 + \dots + x_n$ . Show that Alice's posterior distribution is  $\text{Beta}(\alpha + s_n, \beta + n - s_n)$ . Namely,

$$\Theta | \mathcal{D}_n \sim \text{Beta}(\alpha + s_n, \beta + n - s_n). \quad (466)$$

**Exercise A.7.8.** Let  $Y \sim \text{Uniform}([0, 1])$  and  $X \sim \text{Binomial}(n, Y)$  be independent RVs.

(i) Use iterated expectation for probability to write

$$\mathbb{P}(X = k) = \binom{n}{k} \int_0^1 y^k (1-y)^{n-k} dy. \quad (467)$$

(ii) Write  $A_{n,k} = \int_0^1 y^k (1-y)^{n-k} dy$ . Use integration by parts and show that

$$A_{n,k} = \frac{k}{n-k+1} A_{n,k-1}. \quad (468)$$

for all  $1 \leq k \leq n$ . Conclude that for all  $0 \leq k \leq n$ ,

$$A_{n,k} = \frac{1}{\binom{n}{k}} \frac{1}{n+1}. \quad (469)$$

(iii) Conclude that  $X \sim \text{Uniform}(\{0, 1, \dots, n\})$ .

**Exercise A.7.9** (Beta distribution). A random variable  $X$  taking values from  $[0, 1]$  has Beta distribution of parameters  $\alpha$  and  $\beta$ , which we denote by  $\text{Beta}(\alpha, \beta)$ , if it has PDF

$$f_X(x) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}, \quad (470)$$

where  $\Gamma(z)$  is the Euler Gamma function defined by

$$\Gamma(z) = \int_0^\infty x^{z-1} e^{-x} dx. \quad (471)$$

(i) Use integration by parts to show the following recursion

$$\Gamma(z+1) = z\Gamma(z). \quad (472)$$

Deduce that  $\Gamma(n) = (n-1)!$  for all integers  $n \geq 1$ .

(ii) Let  $X \sim \text{Beta}(k+1, n-k+1)$ . Use (i) to show that

$$f_X(x) = \frac{n!(n+1)}{k!(n-k)!} x^k (1-x)^{n-k} = \frac{x^k (1-x)^{n-k}}{1/\binom{n}{k} (n+1)}. \quad (473)$$

Use Exercise A.7.8 to verify that the above function is indeed a PDF (i.e., it integrates to 1).

(iii)\* Show that if  $X \sim \text{Beta}(\alpha, \beta)$ , then

$$\mathbb{E}[X] = \frac{\alpha}{\alpha + \beta}, \quad \text{Var}(X) = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}. \quad (474)$$

**Example A.7.10** (Bayesian estimator). Recall Exercise A.7.7, where we want to infer an unknown parameter  $\theta$  in  $\text{Bernoulli}(\theta)$ , starting from the beta prior  $\pi = \text{Beta}(\alpha, \beta)$ . After observing a data  $\mathcal{D} = (x_1, \dots, x_n)$  consisting of  $s_n = \sum_{i=1}^n x_i$  successes, our posterior distribution is Beta distribution with parameters  $s_n + 1$  and  $\beta + n - s_n$ . In other words,

$$\Theta | \mathcal{D} \sim \text{Beta}(\alpha + s_n, \beta + n - s_n). \quad (475)$$

Now suppose we want to get a point estimator  $\hat{\theta}$  for the unknown parameter from our posterior distribution. Any choice of the estimator will result in some kind of error, so we would like our estimator to be minimizing an error function of choice. A standard choice is the mean squared error (MSE), which in our case would be under conditioning on the observed data. That is, we want our estimator  $\hat{\theta}$  to be such that

$$\hat{\theta} = \operatorname{argmin}_{\theta} \mathbb{E}[(\Theta - \hat{\theta})^2 | \mathcal{D}]. \quad (476)$$

According to Exercise A.4.2, the above MSE will be minimized when  $\hat{\theta} = \mathbb{E}[\Theta | \mathcal{D}]$ , and the minimum MSE is the conditional variance  $\text{Var}(\Theta | \mathcal{D})$ . Noting the mean of beta distribution in Exercise A.7.9, we obtain the following Bayesian estimator

$$\hat{\theta} = \mathbb{E}[\Theta | \mathcal{D}] = \mathbb{E}[\text{Beta}(\alpha + s_n, \beta + n - s_n)] = \frac{\alpha + s_n}{\alpha + \beta + n}. \quad (477)$$

