

# Day 1

## Classification

This day will serve as an introduction to machine learning. We recall some fundamental concepts about decision theory and classification. We also present some widely used models and algorithms and try to provide the main motivation behind them. There are several textbooks that provide a thorough description of some of the concepts introduced here: for example, Mitchell (1997), Duda et al. (2001), Schölkopf and Smola (2002), Joachims (2002), Bishop (2006), Manning et al. (2008), to name just a few. The concepts that we introduce in this chapter will be revisited in later chapters, where the same algorithms and models will be adapted to structured inputs and outputs. For now, we concern only with multi-class classification (with just a few classes).

### 1.1 Notation

In what follows, we denote by  $\mathcal{X}$  our *input set* (also called *observation set*), and by  $\mathcal{Y}$  our *output set*. We will make no assumptions about the set  $\mathcal{X}$ , which can be continuous or discrete. In this lecture, we consider *classification* problems, where  $\mathcal{Y} = \{c_1, \dots, c_K\}$  is a finite set, consisting of  $K$  *classes* (also called *labels*). For example,  $\mathcal{X}$  can be a set of documents in natural language, and  $\mathcal{Y}$  a set of topics, the goal being to assign a topic to each document.

We use upper-case letters for denoting random variables, and lower-case letters for value assignments to those variables: for example,

- $X$  is a random variable taking values on  $\mathcal{X}$ ,
- $Y$  is a random variable taking values on  $\mathcal{Y}$ ,
- $x \in \mathcal{X}$  and  $y \in \mathcal{Y}$  are particular values for  $X$  and  $Y$ .

We consider *events* such as  $X = x$ ,  $Y = y$ , etc. Throughout, we use modified notation and let  $P(y)$  denote the *probability* associated with the event  $Y = y$

(instead of writing  $P_Y(Y = y)$ ). *Joint* and *conditional* probabilities are denoted respectively as  $P(x, y) \triangleq P_{X,Y}(X = x \wedge Y = y)$  and  $P(x|y) \triangleq P_{X|Y}(X = x | Y = y)$ . From the laws of probabilities:

$$P(x, y) = P(y|x)P(x), \quad (1.1)$$

for all  $x \in \mathcal{X}$  and  $y \in \mathcal{Y}$ .

Quantities that are predicted or estimated from the data will be appended a hat-symbol: for example, estimations of the probabilities above are denoted as  $\hat{P}(y)$ ,  $\hat{P}(x, y)$  and  $\hat{P}(y|x)$ ; and a prediction of an output will be denoted  $\hat{y}$ .

We assume that a *training dataset*  $\mathcal{D}$  is provided which consists of input-output pairs (called *examples* or *instances*):

$$\mathcal{D} = \{(x^1, y^1), \dots, (x^M, y^M)\} \subseteq \mathcal{X} \times \mathcal{Y}. \quad (1.2)$$

The goal of (supervised) machine learning is to use  $\mathcal{D}$  to learn a function  $h$  (called a *classifier*) that maps from  $\mathcal{X}$  to  $\mathcal{Y}$ : this way, given a new instance  $x \in \mathcal{X}$  (test example), the machine makes a prediction  $\hat{y}$  by evaluating  $h$  on  $x$ , i.e.,  $\hat{y} = h(x)$ .

## 1.2 Generative Classifiers: Naïve Bayes

If we knew the *true* distribution  $P(X, Y)$ , the best possible classifier (Bayes optimal) would be one which predicts according to

$$\begin{aligned} \hat{y} &= \arg \max_{y \in \mathcal{Y}} P(y|x) = \arg \max_{y \in \mathcal{Y}} \frac{P(x, y)}{P(x)} \\ &\stackrel{\dagger}{=} \arg \max_{y \in \mathcal{Y}} P(x, y) \\ &= \arg \max_{y \in \mathcal{Y}} P(y)P(x|y), \end{aligned} \quad (1.3)$$

where in  $\dagger$  we used the fact that  $P(x)$  is constant with respect to  $y$ . The probability distributions  $P(Y)$  and  $P(X|Y)$  are respectively called the *class prior* and the *class conditionals*. Figure 1.2 shows an example of the Bayes optimal decision boundary for a toy example. Generative models assume data are generated according to the following generative story (independently for each  $m = 1, \dots, M$ ):

1. A class  $y_m \sim P(Y)$  is drawn from the class prior distribution;
2. An input  $x_m \sim P(X|Y = y_m)$  is drawn from the corresponding class conditional.

Training a generative model amounts to *estimating* these probabilities using the dataset  $\mathcal{D}$ , yielding estimates  $\hat{P}(y)$  and  $\hat{P}(x|y)$ .

Simple Data Set -- Mean1= (-1.00,-1.00) Var1 = 1.00 Mean2= (1.00,1.00) Var2= 1.00  
Nr. Points=100.00, Balance=0.50 Train-Dev-Test (0.80,.00,0.20)

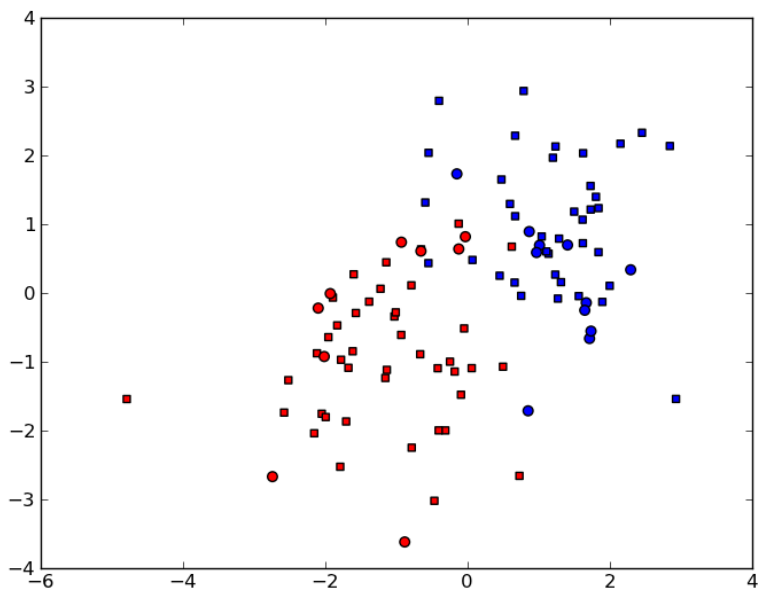


Figure 1.1: Example of a dataset. The input set consists in points in the real plane,  $\mathcal{X} = \mathbb{R}^2$ , and the output set consists of two classes (Red and Blue). Training points are represented as squares, while test points are represented as circles.

Simple Data Set -- Mean1= (-1.00,-1.00) Var1 = 0.50 Mean2= (1.00,1.00) Var2= 0.50  
Nr. Points=100.00, Balance=0.50 Train-Dev-Test (0.80,,0.00,0.20)

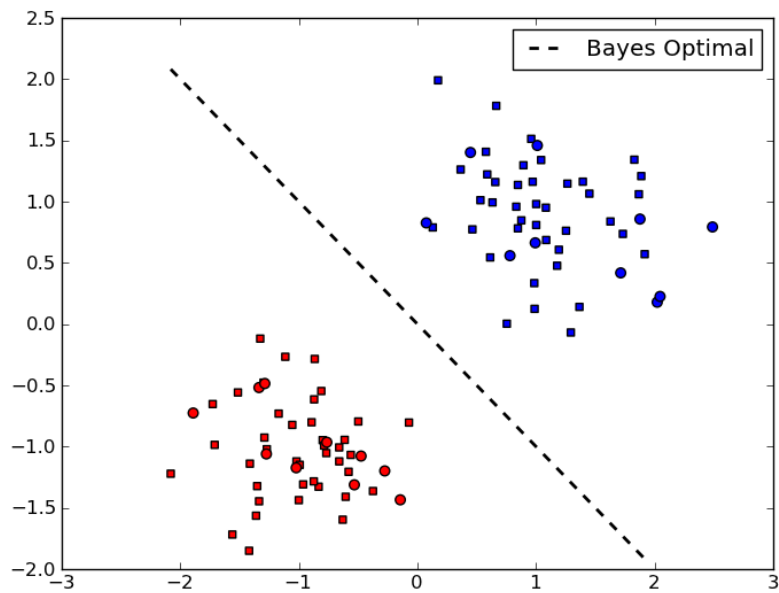


Figure 1.2: Example of a dataset together with the corresponding Bayes optimal decision boundary. The input set consists in points in the real plane,  $\mathcal{X} = \mathcal{R}$ , and the output set consists of two classes (Red and Blue). Training points are represented as squares, while test points are represented as circles.

At test time, given a new input  $x \in \mathcal{X}$ , a prediction is made according to

$$\hat{y} = \arg \max_{y \in \mathcal{Y}} \hat{P}(y) \hat{P}(x|y). \quad (1.4)$$

We are left with two important problems:

1. How should the distributions  $\hat{P}(Y)$  and  $\hat{P}(X|Y)$  be “defined” (i.e., what kind of independence assumptions should they state, or how should they factor?)
2. How should parameters be estimated from the training data  $\mathcal{D}$ ?

The first problem strongly depends on the application at hand. Quite often, there is a natural decomposition of the input variable  $X$  into  $J$  components,

$$X = (X_1, \dots, X_J). \quad (1.5)$$

The naïve Bayes method makes the following assumption:  $X_1, \dots, X_J$  are *conditionally independent given the class*. Mathematically, this means that

$$P(X|Y) = \prod_{j=1}^J P(X_j|Y). \quad (1.6)$$

Note that this independence assumption greatly reduces the number of parameters to be estimated (degrees of freedom) from  $O(\exp(J))$  to  $O(J)$ , hence estimation of  $\hat{P}(Y)$  and  $\hat{P}(X|Y)$  becomes much simpler, as we shall see. It also makes the overall computation much more efficient (in particular for large  $J$ ) and it decreases the risk of overfitting the data. On the other hand, if the assumption is over-simplistic it may increase the risk of under-fitting.

The *maximum likelihood criterion* aims to maximize the probability of the training sample, assuming it was generated iid. This probability (call it  $P(\mathcal{D})$ ) factorizes as

$$\begin{aligned} P(\mathcal{D}) &= \prod_{m=1}^M P(x^m, y^m) \\ &= \prod_{m=1}^M P(y^m) \prod_{j=1}^J P(x_j^m | y^m). \end{aligned} \quad (1.7)$$

### 1.2.1 Example: 2-D Gaussians

We first illustrate the naïve Bayes assumption with a toy example. Suppose that  $\mathcal{X} = \mathbb{R}^2$  and  $\mathcal{Y} = \{1, 2\}$ . Assume that each class-conditional is a two-dimensional Gaussian distribution with fixed covariance, i.e.,  $P(X_1, X_2|Y = y) = \mathcal{N}(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y)$ .

According to the naïve Bayes assumption,  $\hat{P}(X_1, X_2|Y) = \hat{P}(X_1|Y)\hat{P}(X_2|Y)$  (remark: this is equivalent to assuming that the  $\Sigma_y$  are diagonal!). For simplicity, we also assume that the two classes have unit variance. Then, we have  $\hat{P}(X_1|Y = y) = \mathcal{N}(\mu_{y1}, 1.0)$  and  $\hat{P}(X_2|Y = y) = \mathcal{N}(\mu_{y2}, 1.0)$  (Figure 1.1 shows an example a dataset of two gaussians with unit variance. Where  $\mu_{y1} = [-1, -1]$  and  $\mu_{y1} = [1, 1]$ . Figure 1.2 shows the same example but where both gaussian have  $\Sigma = 0.5$ , together with the Bayes optimal decision boundary). The parameters that need to be estimated are the class-conditional means  $\mu_{11}, \mu_{12}, \mu_{21}, \mu_{22}$  and the class priors  $\hat{P}(Y = 1)$  and  $\hat{P}(Y = 2)$ . Given a training sample  $\mathcal{D} = \{(x^1, y^1), \dots, (x^M, y^M)\}$ , denote by  $\mathcal{J}_1 \subseteq \{1, \dots, M\}$  the indices of those instances belonging to class 1, and by  $\mathcal{J}_2 \subseteq \{1, \dots, M\}$  the indices of the ones that belong to class 2. The maximum likelihood estimates of the quantities above are:

$$\begin{aligned} \hat{P}(Y = 1) &= \frac{|\mathcal{J}_1|}{M}, & \hat{P}(Y = 2) &= \frac{|\mathcal{J}_2|}{M} \\ \mu_{11} &= \frac{1}{|\mathcal{J}_1|} \sum_{m \in \mathcal{J}_1} x_1^m, & \mu_{12} &= \frac{1}{|\mathcal{J}_1|} \sum_{m \in \mathcal{J}_1} x_2^m \\ \mu_{21} &= \frac{1}{|\mathcal{J}_2|} \sum_{m \in \mathcal{J}_2} x_1^m, & \mu_{22} &= \frac{1}{|\mathcal{J}_2|} \sum_{m \in \mathcal{J}_2} x_2^m. \end{aligned} \quad (1.8)$$

In words: the class priors' estimates are their relative frequencies, and the class-conditional means' estimates are the sample means.

**Exercise 1.1** Start by importing all the libraries necessary for this lab through the following preamble:

```
1 import sys
  sys.path.append("readers/" )
3 sys.path.append("classifiers/" )

5 import simple_data_set as sds
  import linear_classifier as lcc
7 import gaussian_naive_bayes as gnbc
  import naive_bayes as nb
```

Now, generate a training and a test dataset like in the previous example, each with  $M = 100$  points, 50 of each class. Assume the following class-conditionals:  $P(X|Y = 1) \sim \mathcal{N}((-1, -1), \sigma^2 \mathbf{I})$  and  $P(X|Y = 2) \sim \mathcal{N}((1, 1), \mathbf{I})$ , for  $\sigma = 1.0$ . To do this, run the following command from the code directory:

```
sd = sds.SimpleDataSet(nr_examples=100, g1 = [[-1, -1], 1],
  g2 = [[1, 1], 1], balance=0.5, split=[0.5, 0, 0.5])
```

You can visualize your data and see the Bayes optimal surface boundary by typing:

```
1 fig,axis = sd.plot_data()
```

Now, run naïve Bayes on this dataset. To do that, use the class `GaussianNaiveBayes`, which is defined in the file `GaussianNaiveBayes.py` under the `classification` directory. Report your estimates, as well as training set and testing set accuracies:

```
1 gnb = gnb.GaussianNaiveBayes()
2 params_nb_sd = gnb.train(sd.train_X, sd.train_y)
3
4 print "Estimated Means"
5 print gnb.means
6 print "Estimated Priors"
7 print gnb.prior
8 y_pred_train = gnb.test(sd.train_X, params_nb_sd)
9 acc_train = gnb.evaluate(sd.train_y, y_pred_train)
10 y_pred_test = gnb.test(sd.test_X, params_nb_sd)
11 acc_test = gnb.evaluate(sd.test_y, y_pred_test)
12 print "Gaussian Naive Bayes Simple Dataset Accuracy train:
13     %f test: %f"%(acc_train, acc_test)
```

To visualize the surface boundary estimated by naïve Bayes, type:

```
fig,axis = sd.add_line(fig,axis,params_nb_sd,"Naive Bayes",
"red")
```

Do not worry for now about why the surface boundaries look the way they look. This is going to be the subject of §1.3.

Repeat the exercise above for different values of  $\sigma^2$ , different balances, and different sample sizes. What do you observe?

## 1.2.2 Example: Multinomial Model for Document Classification

We now consider a more realistic scenario where the naïve Bayes classifier may be applied. Suppose that the task is *document classification*:  $\mathcal{X}$  is the set of all possible documents, and  $\mathcal{Y} = \{c_1, \dots, c_K\}$  is a set of *topics* for those documents. Let  $\mathcal{V} = \{w_1, \dots, w_J\}$  be the vocabulary, i.e., the set of words that occur in some document.

A very popular document representation is through a “bag-of-words”: each document is seen as a multiset of words along with their frequencies; word ordering is ignored. We are going to see that this is equivalent to a naïve Bayes assumption with the *multinomial model*.<sup>1</sup> We associate to each class a multinomial distribution, which ignores word ordering, but takes into consideration the frequency with which each word appears in a document. For simplicity, we assume that all documents have the same length  $L$ .<sup>2</sup> Each document  $x$  is assumed to have been generated as follows. First, a class  $y$  is generated according to  $P(y)$ . Then,  $x$  is generated by sequentially picking words from  $\mathcal{V}$  with replacement. Each word  $w_j$  is picked with probability  $P(w_j|y)$ . For example, the probability of generating a document  $x = w_{j_1} \dots w_{j_L}$  (i.e., a sequence of  $L$  words—*tokens*— $w_{j_1}, \dots, w_{j_L}$ ) is

$$P(x|y) = \prod_{l=1}^L P(w_{j_l}|y) = \prod_{j=1}^J P(w_j|y)^{n_j(x)}, \quad (1.9)$$

where  $n_j(x)$  is the number of occurrences of word  $w_j$  in document  $x$ .

Hence, the assumption is that word occurrences (*tokens*) are independent given the class. The parameters that need to be estimated are  $\hat{P}(c_1), \dots, \hat{P}(c_K)$ , and  $\hat{P}(w_j|c_k)$  for  $j = 1, \dots, J$  and  $k = 1, \dots, K$ . Given a training sample  $\mathcal{D} = \{(x^1, y^1), \dots, (x^M, y^M)\}$ , denote by  $\mathcal{J}_k$  the indices of those instances belonging to the  $k$ th class. The maximum likelihood estimates of the quantities above are:

$$\hat{P}(c_k) = \frac{|\mathcal{J}_k|}{M}, \quad \hat{P}(w_j|c_k) = \frac{\sum_{m \in \mathcal{J}_k} n_j(x^m)}{|\mathcal{J}_k|}. \quad (1.10)$$

In words: the class priors’ estimates are their relative frequencies (as before), and the class-conditional word probabilities are the relative frequencies of those words across documents with that class.

**Exercise 1.2** *In this exercise we will use the the Amazon sentiment analysis data (Blitzer et al., 2007), where the goal is to classify text documents as expressing a positive or negative sentiment (i.e., a classification problem with two labels). We are going to focus on book reviews. To load the data, type:*

```
1 import sentiment_reader as srs
2 import naive_bayes as nb
3
4 scr = srs.SentimentCorpus("books")
```

<sup>1</sup>Another popular model for documents is the Bernoulli model, which only looks at the presence/absence of a word in a document, rather than word frequency. See Manning et al. (2008); McCallum and Nigam (1998) for further information.

<sup>2</sup>We can get rid of this assumption by defining a distribution on the document length. Everything stays the same if that distribution is uniform up to a maximum document length.



This will load the data in a bag-of-words representation where rare words (occurring less than 5 times in the training data) are removed.

1. Create a file `MultinomialNaiveBayes.py` and implement the naïve Bayes with the multinomial model in a new class `MultinomialNaiveBayes`. (Hint: look at the implementation of `GaussianNaiveBayes` for inspiration).
2. Run naïve Bayes with the multinomial model on the Amazon dataset (sentiment classification) and report results both for training and testing:

```
import multinomial_naive_bayes as mnb
2
mnb = mnb.MultinomialNaiveBayes()
4 params_nb_sc = mnb.train(scr.train_X,scr.train_y)
y_pred_train = mnb.test(scr.train_X,params_nb_sc)
6 acc_train = mnb.evaluate(scr.train_y, y_pred_train)
y_pred_test = mnb.test(scr.test_X,params_nb_sc)
8 acc_test = mnb.evaluate(scr.test_y, y_pred_test)
print "Multinomial Naive Bayes Amazon Sentiment Accuracy
train: %f test: %f"%(acc_train,acc_test)
```

3. Observe that words that were not observed at training time cause problems at test time. Why? To solve this problem, apply a simple add-one smoothing technique: replace the expression in Eq. 1.10 for the estimation of the conditional probabilities by

$$\hat{P}(w_j|c_k) = \frac{1 + \sum_{m \in \mathcal{J}_k} n_j(x^m)}{J + |\mathcal{J}_k|}.$$

where  $J$  is the number of distinct words.

This is a widely used smoothing strategy which has a Bayesian interpretation: it corresponds to choosing a uniform prior for the word distribution on both classes, and to replace the maximum likelihood criterion by a maximum a posteriori approach. This is a form of regularization, preventing the model from overfitting on the training data. See e.g. Manning and Schütze (1999); Manning et al. (2008) for more information. Report the new accuracies.

### 1.3 Features and Linear Classifiers

In the previous section, we assumed a particular representation for the input objects  $x \in \mathcal{X}$ : points in a 2D Euclidean space (in the Gaussian example) and bag-of-words representations of text documents (in the sentiment data example).

The methods discussed in this lecture are also applicable to a wide range of problems, regardless of the intricacies of our input objects. It is useful to think

about each  $x \in \mathcal{X}$  as an abstract object, which is subject to a set of descriptions or measurements, which are called *features*. A feature is simply a real number that describes the value of some property of  $x$ . For instance in the toy examples described above you can think of  $\mathcal{X}$  as a set of points and the features to be its 2D coordinates. Let  $g_1(x), \dots, g_J(x)$  be  $J$  features of  $x$ . We call the vector

$$\mathbf{g}(x) = (g_1(x), \dots, g_J(x)) \quad (1.11)$$

a *feature vector representation* of  $x$ . The map  $\mathbf{g} : \mathcal{X} \rightarrow \mathbb{R}^J$  is called a *feature mapping*.

In NLP applications, features are often binary-valued and result from evaluating propositions such as:

$$g_1(x) \triangleq \begin{cases} 1, & \text{if sentence } x \text{ contains the word } \textit{Ronaldo} \\ 0, & \text{otherwise.} \end{cases} \quad (1.12)$$

$$g_2(x) \triangleq \begin{cases} 1, & \text{if all words in sentence } x \text{ are capitalized} \\ 0, & \text{otherwise.} \end{cases} \quad (1.13)$$

$$g_3(x) \triangleq \begin{cases} 1, & \text{if } x \text{ contains any of the words } \textit{amazing}, \textit{excellent} \text{ or } \textit{:} \\ 0, & \text{otherwise.} \end{cases} \quad (1.14)$$

In this example, the feature vector representation of the sentence  $x$

Ronaldo kicked the ball and scored an amazing goal!

would be  $\mathbf{g}(x) = (1, 0, 1)$ .

In multi-class learning problems, rather than associating features only with the input objects, it is useful to consider *joint feature mappings*  $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^D$ . In that case, the *joint feature vector*  $f(x, y)$  can be seen as a collection of joint input-output measurements. For example:

$$f_1(x, y) \triangleq \begin{cases} 1, & \text{if } x \text{ contains } \textit{Ronaldo}, \text{ and topic } y \text{ is } \textit{sport} \\ 0, & \text{otherwise.} \end{cases} \quad (1.15)$$

$$f_2(x, y) \triangleq \begin{cases} 1, & \text{if } x \text{ contains } \textit{Ronaldo}, \text{ and topic } y \text{ is } \textit{politics} \\ 0, & \text{otherwise.} \end{cases} \quad (1.16)$$

A very simple form of defining a joint feature mapping which is often employed is via:

$$\begin{aligned} f(x, y) &\triangleq \mathbf{g}(x) \otimes \mathbf{e}_y \\ &= (0, \dots, 0, \underbrace{\mathbf{g}(x)}_{y\text{th slot}}, 0, \dots, 0) \end{aligned} \quad (1.17)$$

where  $\mathbf{g}(x) \in \mathbb{R}^J$  is an input feature vector,  $\otimes$  is the Kronecker product ( $[\mathbf{a} \otimes \mathbf{b}]_{ij} = a_i b_j$ ) and  $\mathbf{e}_y \in \mathbb{R}^K$ , with  $[e_y]_c = 1$  iff  $y = c$ , and 0 otherwise. Hence

$f(x, y) \in \mathbb{R}^D$  with  $D = JK$ .

Linear classifiers are very popular in natural language processing applications. They make their decision based on the rule:

$$\hat{y} = \arg \max_{y \in \mathcal{Y}} \boldsymbol{w} \cdot \boldsymbol{f}(x, y). \quad (1.18)$$

where

- $\boldsymbol{w} \in \mathbb{R}^D$  is a *weight vector*;
- $\boldsymbol{f}(x, y) \in \mathbb{R}^D$  is a *feature vector*;
- $\boldsymbol{w} \cdot \boldsymbol{f}(x, y) = \sum_{d=1}^D w_d f_d(x, y)$  is the inner product between  $\boldsymbol{w}$  and  $\boldsymbol{f}(x, y)$ .

Hence, each feature  $f_d(x, y)$  has a weight  $w_d$  and, for each class  $y \in \mathcal{Y}$ , a score is computed by linearly combining all the weighted features. All these scores are compared, and a prediction is made by choosing the class with the largest score.

**Remark 1.1** With the design above (Eq. 1.17), and decomposing the weight vector as  $\boldsymbol{w} = (\boldsymbol{w}_{c_1}, \dots, \boldsymbol{w}_{c_K})$ , we have that

$$\boldsymbol{w} \cdot \boldsymbol{f}(x, y) = \boldsymbol{w}_y \cdot \boldsymbol{g}(x). \quad (1.19)$$

In words: each class  $y \in \mathcal{Y}$  gets its own weight vector  $\boldsymbol{w}_y$ , and one defines an input feature vector  $\boldsymbol{g}(x)$  that only looks at the input  $x \in \mathcal{X}$ . This representation is very useful when features only depend on input  $x$  since it allows a more compact representation. Note that the number of features is normally very large.

**Remark 1.2** The multinomial naïve Bayes classifier described in the previous section is an instance of a linear classifier (in fact, so is the 2-D Gaussian Bayes classifier—try to show this). Recall that the naïve Bayes classifier predicts according to  $\hat{y} = \arg \max_{y \in \mathcal{Y}} \hat{P}(y) \hat{P}(x|y)$ . Taking logs, in the multinomial model for document classification this is equivalent to:

$$\begin{aligned} \hat{y} &= \arg \max_{y \in \mathcal{Y}} \log \hat{P}(y) + \log \hat{P}(x|y) \\ &= \arg \max_{y \in \mathcal{Y}} \log \hat{P}(y) + \sum_{j=1}^J n_j(x) \log \hat{P}(w_j|y) \\ &= \arg \max_{y \in \mathcal{Y}} \boldsymbol{w}_y \cdot \boldsymbol{g}(x), \end{aligned} \quad (1.20)$$

where

$$\begin{aligned} \boldsymbol{w}_y &= (b_y, \log \hat{P}(w_1|y), \dots, \log \hat{P}(w_J|y)) \\ b_y &= \log \hat{P}(y) \\ \boldsymbol{g}(x) &= (1, n_1(x), \dots, n_J(x)). \end{aligned} \quad (1.21)$$

---

**Algorithm 2** Averaged perceptron

---

- 1: **input:** dataset  $\mathcal{D}$ , number of rounds  $R$
- 2: initialize  $t = 0, \mathbf{w}^t = \mathbf{0}$
- 3: **for**  $r = 1$  **to**  $R$  **do**
- 4:    $\mathcal{D}_s = \text{shuffle}(\mathcal{D})$
- 5:   **for**  $i = 1$  **to**  $M$  **do**
- 6:      $m = \mathcal{D}_s(i)$
- 7:      $t = t + 1$
- 8:     take training pair  $(x^m, y^m)$  and predict using the current model:

$$\hat{y} \leftarrow \arg \max_{y' \in \mathcal{Y}} \mathbf{w}^t \cdot \mathbf{f}(x^m, y')$$

- 9:     update the model:  $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t + \mathbf{f}(x^m, y^m) - \mathbf{f}(x^m, \hat{y})$
  - 10:   **end for**
  - 11: **end for**
  - 12: **output:** the averaged model  $\hat{\mathbf{w}} \leftarrow \frac{1}{t} \sum_{i=1}^t \mathbf{w}^i$
- 

Hence, the multinomial model yields a prediction rule of the form

$$\hat{y} = \arg \max_{y \in \mathcal{Y}} \mathbf{w}_y \cdot \mathbf{g}(x). \quad (1.22)$$

**Exercise 1.3** Show that the Gaussian naïve Bayes classifier with shared and given variance is also a linear classifier, and derive the formulas for  $\mathbf{w}_y, b_y$ . You should obtain the formulas that are implemented in the `train` method of `GaussianNaiveBayes`.

Look again at the decision boundary that you have found in Exercise 1.1 and compare it with the Bayes optimal classifier.

## 1.4 Online Algorithms: Perceptron and MIRA

### 1.4.1 Perceptron

Perhaps the oldest algorithm to train a linear classifier is the *perceptron* (Rosenblatt, 1958), which we depict as Alg. 2.<sup>3</sup>

The perceptron algorithm works as follows: at each round, it takes an input datum, and uses the current model to make a prediction. If the prediction is correct, nothing happens. Otherwise, the model is corrected by adding the feature vector w.r.t. the correct output and subtracting the feature vector w.r.t. the predicted (wrong) output. Then, we proceed to the next round. Alg. 2 is remarkably simple; yet it often reaches a very good performance, often bet-

---

<sup>3</sup>Actually, we are showing a more robust variant of the perceptron, which averages the weight vector as a post-processing step.

ter than the Naïve Bayes model, and usually not much worse than maximum entropy models or SVMs (which will be described in the next section).

A weight vector  $w$  defines a *separating hyperplane* if it classifies all the training data correctly, *i.e.*, if  $y^m = \arg \max_{y \in \mathcal{Y}} w \cdot f(x^m, y)$  hold for  $m = 1, \dots, M$ . A dataset  $\mathcal{D}$  is *separable* if such a weight vector exists (in general,  $w$  is not unique). A very important property of the perceptron algorithm is the following: if  $\mathcal{D}$  is separable, then the number of mistakes made by the perceptron algorithm is *finite*. This means that under this assumption, the perceptron will eventually reach a separating hyperplane  $w$ .

There are other variants of the perceptron (e.g., with regularization) which we omit for brevity.

**Exercise 1.4** We provide an implementation of the perceptron algorithm in the class `Perceptron` (file `Perceptron.py`).

1. Run the perceptron algorithm on the simple dataset previously generated and report its train and test set accuracy:

```
1 import perceptron as percc
3 perc = percc.Perceptron()
  params_perc_sd = perc.train(sd.train_X, sd.train_y)
5 y_pred_train = perc.test(sd.train_X, params_perc_sd)
  acc_train = perc.evaluate(sd.train_y, y_pred_train)
7 y_pred_test = perc.test(sd.test_X, params_perc_sd)
  acc_test = perc.evaluate(sd.test_y, y_pred_test)
9 print "Perceptron Simple Dataset Accuracy train: %f test: %f" %
      (acc_train, acc_test)
```

2. Plot the decision boundary found:

```
1 fig, axis = sd.add_line(fig, axis, params_perc_sd, "Perceptron",
  , "blue")
```

Change the code to save the intermediate weight vectors, and plot them every five iterations. What do you observe?

3. Run the perceptron algorithm on the Amazon dataset.

## 1.4.2 Margin Infused Relaxed Algorithm (MIRA)

The MIRA algorithm (Crammer and Singer, 2002; Crammer et al., 2006) has achieved very good performance in NLP problems. At each round  $t$ , MIRA updates the weight vector by solving the following optimization problem:

$$\mathbf{w}^{t+1} \leftarrow \arg \min_{\mathbf{w}, \xi} \quad \xi + \frac{\lambda}{2} \|\mathbf{w} - \mathbf{w}^t\|^2 \quad (1.23)$$

$$\text{s.t. } \mathbf{w} \cdot \mathbf{f}(x^m, y^m) \geq \mathbf{w} \cdot \mathbf{f}(x^m, \hat{y}) + 1 - \xi \quad (1.24)$$

$$\xi \geq 0, \quad (1.25)$$

where  $\hat{y} = \arg \max_{y' \in \mathcal{Y}} \mathbf{w}^t \cdot \mathbf{f}(x^m, y')$  is the prediction using the model with weight vector  $\mathbf{w}^t$ . By inspecting Eq. 1.23 we see that MIRA attempts to achieve a tradeoff between *conservativeness* (penalizing large changes from the previous weight vector via the term  $\frac{\lambda}{2} \|\mathbf{w} - \mathbf{w}^t\|^2$ ) and *correctness* (by requiring, through the constraints, that the new model  $\mathbf{w}^{t+1}$  “separates” the true output from the prediction with a margin (although slack  $\xi \geq 0$  is allowed)).<sup>4</sup> Note that, if the prediction is correct ( $\hat{y} = y^m$ ) the solution of the problem Eq. 1.23 leaves the weight vector unchanged ( $\mathbf{w}^{t+1} = \mathbf{w}^t$ ). This quadratic programming problem has a closed form solution:<sup>5</sup>

$$\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t + \eta^t (\mathbf{f}(x^m, y^m) - \mathbf{f}(x^m, \hat{y})),$$

with

$$\eta^t = \min \left\{ \lambda^{-1}, \frac{\mathbf{w}^t \cdot \mathbf{f}(x^m, \hat{y}) - \mathbf{w}^t \cdot \mathbf{f}(x^m, y^m) + \rho(\hat{y}, y^m)}{\|\mathbf{f}(x^m, y^m) - \mathbf{f}(x^m, \hat{y})\|^2} \right\},$$

where  $\rho : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$  is a non-negative cost function, such that  $\rho(\hat{y}, y)$  is the cost incurred by predicting  $\hat{y}$  when the true output is  $y$ ; we assume  $\rho(y, y) = 0$  for all  $y \in \mathcal{Y}$ . For simplicity, we focus here on the 0/1-cost (but keep in mind that other cost functions are possible):

$$\rho(\hat{y}, y) = \begin{cases} 1 & \text{if } \hat{y} \neq y \\ 0 & \text{otherwise.} \end{cases} \quad (1.26)$$

MIRA is depicted in Alg. 3. For other variants of MIRA, see Crammer et al. (2006).

**Exercise 1.5** Implement the MIRA algorithm (Hint: use the perceptron algorithm as a starting point and modify it as necessary). Do this by creating a file `Mira.py` and implement class `Mira`. Then, repeat the perceptron exercise now using MIRA, for several values of  $\lambda$ :

```
1 import mira as mirac
3 mira = mirac.Mira()
  mira.regularizer = 1.0 # This is lambda
```

<sup>4</sup>The intuition for this large margin separation is the same for support vector machines, which will be discussed in §1.5.2.

<sup>5</sup>Note that the perceptron updates are identical, except that we always have  $\eta_t = 1$ .

---

**Algorithm 3** MIRA

---

- 1: **input:** dataset  $\mathcal{D}$ , parameter  $\lambda$ , number of rounds  $R$
- 2: initialize  $t = 0, \mathbf{w}^t = \mathbf{0}$
- 3: **for**  $r = 1$  **to**  $R$  **do**
- 4:    $\mathcal{D}_s = \text{shuffle}(\mathcal{D})$
- 5:   **for**  $i = 1$  **to**  $M$  **do**
- 6:      $m = \mathcal{D}_s(i)$
- 7:      $t = t + 1$
- 8:     take training pair  $(x^m, y^m)$  and predict using the current model:

$$\hat{y} \leftarrow \arg \max_{y' \in \mathcal{Y}} \mathbf{w}^t \cdot \mathbf{f}(x^m, y')$$

- 9:     compute loss:  $\ell^t = \mathbf{w}^t \cdot \mathbf{f}(x^m, \hat{y}) - \mathbf{w}^t \cdot \mathbf{f}(x^m, y^m) + \rho(\hat{y}, y^m)$
  - 10:     compute stepsize:  $\eta^t = \min \left\{ \lambda^{-1}, \frac{\ell^t}{\|\mathbf{f}(x^m, y^m) - \mathbf{f}(x^m, \hat{y})\|^2} \right\}$
  - 11:     update the model:  $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t + \eta^t (\mathbf{f}(x^m, y^m) - \mathbf{f}(x^m, \hat{y}))$
  - 12:   **end for**
  - 13: **end for**
  - 14: **output:** the averaged model  $\hat{\mathbf{w}} \leftarrow \frac{1}{t} \sum_{i=1}^t \mathbf{w}^i$
- 

```
5 params_mira_sd = mira.train(sd.train_X, sd.train_y)
  y_pred_train = mira.test(sd.train_X, params_mira_sd)
7 acc_train = mira.evaluate(sd.train_y, y_pred_train)
  y_pred_test = mira.test(sd.test_X, params_mira_sd)
9 acc_test = mira.evaluate(sd.test_y, y_pred_test)
  print "Mira Simple Dataset Accuracy train: %f test: %f"%(
    acc_train, acc_test)
11 fig, axis = sd.add_line(fig, axis, params_mira_sd, "Mira", "green")

13 params_mira_sc = mira.train(scr.train_X, scr.train_y)
  y_pred_train = mira.test(scr.train_X, params_mira_sc)
15 acc_train = mira.evaluate(scr.train_y, y_pred_train)
  y_pred_test = mira.test(scr.test_X, params_mira_sc)
17 acc_test = mira.evaluate(scr.test_y, y_pred_test)
  print "Mira Amazon Sentiment Accuracy train: %f test: %f"%(
    acc_train, acc_test)
```

Compare the results achieved and separating hiperplanes found.

## 1.5 Discriminative Classifiers: Maximum Entropy and Support Vector Machines

Unlike the naïve Bayes classifier, the algorithms described in the last section (perceptron and MIRA) directly focus on finding a separating hyperplane to discriminate among the classes, rather than attempting to model the probability  $P(X, Y)$  that generates the data. This kind of methods are called *discriminative* (by opposition to the *generative* ones). This section presents two important discriminative classifiers, with widespread use in NLP applications: maximum entropy and support vector machines.

### 1.5.1 Maximum Entropy Classifiers

The notion of *entropy* in the context of Information Theory (Shannon, 1948) is one of the most significant advances in mathematics in the twentieth century. The principle of *maximum entropy* (which appears under different names, such as “maximum mutual information” or “minimum Kullback-Leibler divergence”) plays a fundamental role in many methods in statistics and machine learning (Jaynes, 1982). For an excellent textbook on Information Theory, we recommend Cover et al. (1991). The basic rationale is that choosing the model with the highest entropy (subject to constraints that depend on the observed data) corresponds to making the fewest possible assumptions regarding what was unobserved, trying to making uncertainty about the model as large as possible. For example, if we have a dice and want to estimate the probability of its outcomes, the distribution with the highest entropy would be the uniform distribution (each outcome having of probability a  $1/6$ ). Now suppose that we partition the set of possible outcomes in two groups and are only told about how many times outcomes on each of the groups have occurred. If we know that outcomes  $\{1, 2, 3\}$  occurred 10 times in total, and  $\{4, 5, 6\}$  occurred 30 times in total, then the principle of maximum entropy would lead us to estimate  $P(1) = P(2) = P(3) = 1/12$  and  $P(4) = P(5) = P(6) = 1/4$  (i.e., outcomes would be uniform within each of the two groups).

For an introduction of maximum entropy models, along with pointers to the literature, see <http://www.cs.cmu.edu/~abberger/maxent.html>. A fundamental result is that the maximum entropy distribution  $P_w(Y|X)$  under *first moment matching constraints* (which mean that feature expectations under that distribution  $\frac{1}{M} \sum_m E_{Y \sim P_w} [f(x_m, Y)]$  must match the observed relative frequencies  $\frac{1}{M} \sum_m f(x_m, y_m)$ ) is a *log-linear model*. The dual of that optimization problem is that of maximizing likelihood in a log-linear model (in the binary case, called *logistic regression* model).



The maximum entropy distribution<sup>6</sup> has the following parametric form:

$$P_w(y|x) = \frac{\exp(\mathbf{w} \cdot \mathbf{f}(x, y))}{Z(\mathbf{w}, x)} \quad (1.27)$$

The denominator in Eq. 1.27 is called the *partition function*:

$$Z(\mathbf{w}, x) = \sum_{y' \in \mathcal{Y}} \exp(\mathbf{w} \cdot \mathbf{f}(x, y')). \quad (1.28)$$

An important property of the partition function is that the gradient of its logarithm equals the feature expectations:

$$\begin{aligned} \nabla_{\mathbf{w}} \log Z(\mathbf{w}, x) &= E_w[\mathbf{f}(x, Y)] \\ &= \sum_{y' \in \mathcal{Y}} P_w(y'|x) \mathbf{f}(x, y'). \end{aligned} \quad (1.29)$$

Maximum entropy models are trained *discriminatively*: this means that, instead of maximizing the *joint* likelihood  $P_w(x^1, \dots, x^M, y^1, \dots, y^M)$  (like generative approaches, such as naïve Bayes, do), one maximizes directly the *conditional* likelihood  $P_w(y^1, \dots, y^M | x^1, \dots, x^M)$ . The rationale is that one does not need to worry about modeling the input variables if all we want is an accurate estimate of  $P(Y|X)$ , which is what matters for prediction. The average conditional log-likelihood is:

$$\begin{aligned} \mathcal{L}(\mathbf{w}; \mathcal{D}) &= \frac{1}{M} \log P_w(y^1, \dots, y^M | x^1, \dots, x^M) \\ &= \frac{1}{M} \log \prod_{m=1}^M P_w(y^m | x^m) \\ &= \frac{1}{M} \sum_{m=1}^M \log P_w(y^m | x^m) \\ &= \frac{1}{M} \sum_{m=1}^M (\mathbf{w} \cdot \mathbf{f}(x^m, y^m) - \log Z(\mathbf{w}, x^m)). \end{aligned} \quad (1.30)$$

We try to find the parameters  $\mathbf{w}$  that maximize the log-likelihood  $\mathcal{L}(\mathbf{w}; \mathcal{D})$ ; to avoid overfitting, we add a regularization term that penalizes values of  $\mathbf{w}$  that have a high magnitude. The optimization problem becomes:

$$\begin{aligned} \hat{\mathbf{w}} &= \arg \max_{\mathbf{w}} \mathcal{L}(\mathbf{w}; \mathcal{D}) - \frac{\lambda}{2} \|\mathbf{w}\|^2 \\ &= \arg \min_{\mathbf{w}} -\mathcal{L}(\mathbf{w}; \mathcal{D}) + \frac{\lambda}{2} \|\mathbf{w}\|^2. \end{aligned} \quad (1.31)$$

---

<sup>6</sup>Also called a log-linear model, a Boltzmann distribution, or an exponential family of distributions.

Here we use the squared  $L_2$ -norm as the regularizer,<sup>7</sup> but other norms are possible. The scalar  $\lambda \geq 0$  controls the amount of regularization. Unlike the naïve Bayes examples, this optimization problem does not have a closed form solution in general; hence we need to resort to numerical optimization (see section ??). Let  $F_\lambda(\mathbf{w}; \mathcal{D}) = -\mathcal{L}(\mathbf{w}; \mathcal{D}) + \frac{\lambda}{2} \|\mathbf{w}\|^2$  be the objective function in Eq. 1.31. This function is convex, which implies that a local optimum of Eq. 1.31 is also a global optimum.  $F_\lambda(\mathbf{w}; \mathcal{D})$  is also differentiable: its gradient is

$$\begin{aligned} \nabla_{\mathbf{w}} F_\lambda(\mathbf{w}; \mathcal{D}) &= \frac{1}{M} \sum_{m=1}^M (-f(x^m, y^m) + \nabla_{\mathbf{w}} \log Z(\mathbf{w}, x^m)) + \lambda \mathbf{w} \\ &= \frac{1}{M} \sum_{m=1}^M (-f(x^m, y^m) + E_{\mathbf{w}}[f(x^m, Y)]) + \lambda \mathbf{w}. \end{aligned} \quad (1.32)$$

A batch gradient method to optimize Eq. 1.31 is shown in Alg. 4. Essentially, Alg. 4 iterates through the following updates until convergence:

$$\begin{aligned} \mathbf{w}^{t+1} &\leftarrow \mathbf{w}^t - \eta_t \nabla_{\mathbf{w}} F_\lambda(\mathbf{w}^t; \mathcal{D}) \\ &= (1 - \lambda \eta_t) \mathbf{w}^t + \eta_t \frac{1}{M} \sum_{m=1}^M (f(x^m, y^m) - E_{\mathbf{w}}[f(x^m, Y)]). \end{aligned} \quad (1.33)$$

Convergence is ensured for suitable stepsizes  $\eta_t$ . Monotonic decrease of the objective value can also be ensured if  $\eta_t$  is chosen with a suitable line search method, such as Armijo’s rule (Nocedal and Wright, 1999). In practice, more sophisticated methods exist for optimizing Eq. 1.31, such as conjugate gradient or L-BFGS. The latter is an example of a quasi-Newton method, which only requires gradient information, but use past gradients to try to construct second order (Hessian) approximations.

In large-scale problems (very large  $M$ ) batch methods are slow, *online* or *stochastic* optimization make attractive alternative methods. Stochastic gradient methods make “noisy” gradient updates by considering only a single instance at the time. The resulting algorithm is shown as Alg. 5. At each round  $t$ , an instance  $m(t)$  is chosen, either randomly (stochastic variant) or by cycling through the dataset (online variant). The stepsize sequence must decrease with  $t$ : typically,  $\eta_t = \eta_0 t^{-\alpha}$  for some  $\eta_0 > 0$  and  $\alpha \in [1, 2]$ , tuned in a development partition or with cross-validation.

**Exercise 1.6** We provide an implementation of the L-BFGS algorithm for training maximum entropy models in the class `MaxEnt_batch`, as well as an implementation of the SGD algorithm in the class `MaxEnt_online`.

1. Train a maximum entropy model using L-BFGS on the Simple data set (try

<sup>7</sup>In a Bayesian perspective, this corresponds to choosing independent Gaussian priors  $p(w_d) \sim \mathcal{N}(0; 1/\lambda^2)$  for each dimension of the weight vector.

---

**Algorithm 4** Batch Gradient Descent for Maximum Entropy

---

- 1: **input:**  $\mathcal{D}$ ,  $\lambda$ , number of rounds  $T$ ,  
learning rate sequence  $(\eta_t)_{t=1,\dots,T}$
- 2: initialize  $w^1 = \mathbf{0}$
- 3: **for**  $t = 1$  **to**  $T$  **do**
- 4:   **for**  $m = 1$  **to**  $M$  **do**
- 5:     take training pair  $(x^m, y^m)$  and compute conditional probabilities using the current model, for each  $y' \in \mathcal{Y}$ :

$$P_{w^t}(y'|x^m) = \frac{\exp(w^t \cdot f(x^m, y'))}{Z(w, x)}$$

- 6:     compute the feature vector expectation:

$$E_w[f(x^m, Y)] = \sum_{y' \in \mathcal{Y}} P_{w^t}(y'|x^m) f(x^m, y')$$

- 7:   **end for**
- 8:   choose the stepsize  $\eta_t$  using, e.g., Armijo's rule
- 9:   update the model:

$$w^{t+1} \leftarrow (1 - \lambda\eta_t)w^t + \eta_t M^{-1} \sum_{m=1}^M (f(x^m, y^m) - E_w[f(x^m, Y)])$$

- 10: **end for**
  - 11: **output:**  $\hat{w} \leftarrow w^{T+1}$
- 

different values of  $\lambda$ ). Compare the results with the previous methods. Plot the decision boundary.

```
import max_ent_batch as mebc
2
me_lbfgs = mebc.MaxEnt_batch()
4 me_lbfgs.regularizer = 1.0
params_meb_sd = me_lbfgs.train(sd.train_X, sd.train_y)
6 y_pred_train = me_lbfgs.test(sd.train_X, params_meb_sd)
acc_train = me_lbfgs.evaluate(sd.train_y, y_pred_train)
8 y_pred_test = me_lbfgs.test(sd.test_X, params_meb_sd)
acc_test = me_lbfgs.evaluate(sd.test_y, y_pred_test)
10 print "Max-Ent batch Simple Dataset Accuracy train: %f test
: %f"%(acc_train, acc_test)
12 fig, axis = sd.add_line(fig, axis, params_meb_sd, "Max-Ent-
Batch", "orange")
```

---

**Algorithm 5** SGD for Maximum Entropy

---

- 1: **input:**  $\mathcal{D}$ ,  $\lambda$ , number of rounds  $T$ ,  
learning rate sequence  $(\eta_t)_{t=1,\dots,T}$
- 2: initialize  $w^1 = \mathbf{0}$
- 3: **for**  $t = 1$  **to**  $T$  **do**
- 4: choose  $m = m(t)$  randomly
- 5: take training pair  $(x^m, y^m)$  and compute conditional probabilities using the current model, for each  $y' \in \mathcal{Y}$ :

$$P_{w^t}(y'|x^m) = \frac{\exp(w^t \cdot f(x^m, y'))}{Z(w, x)}$$

- 6: compute the feature vector expectation:

$$E_w[f(x^m, Y)] = \sum_{y' \in \mathcal{Y}} P_{w^t}(y'|x^m) f(x^m, y')$$

- 7: update the model:

$$w^{t+1} \leftarrow (1 - \lambda\eta_t)w^t + \eta_t (f(x^m, y^m) - E_w[f(x^m, Y)])$$

- 8: **end for**
  - 9: **output:**  $\hat{w} \leftarrow w^{T+1}$
- 

2. Train a maximum entropy model using L-BFGS, on the Amazon dataset (try different values of  $\lambda$ ) and report training and test set accuracy. What do you observe?

```
params_meb_sc = me_lbfgs.train(scr.train_X, scr.train_y)
2 y_pred_train = me_lbfgs.test(scr.train_X, params_meb_sc)
acc_train = me_lbfgs.evaluate(scr.train_y, y_pred_train)
4 y_pred_test = me_lbfgs.test(scr.test_X, params_meb_sc)
acc_test = me_lbfgs.evaluate(scr.test_y, y_pred_test)
6 print "Max-Ent Batch Amazon Sentiment Accuracy train: %f
test: %f"%(acc_train, acc_test)
```

3. Now, fix  $\lambda = 1.0$  and train with SGD (you might try to adjust the initial step). Compare the objective values obtained during training with those obtained with L-BFGS. What do you observe?

```
import max_ent_online as meoc
2
me_sgd = meoc.MaxEnt_online()
4 me_sgd.regularizer = 1.0
```

```

params_meo_sc = me_sgd.train(scr.train_X,scr.train_y)
6 y_pred_train = me_sgd.test(scr.train_X,params_meo_sc)
acc_train = me_sgd.evaluate(scr.train_y, y_pred_train)
8 y_pred_test = me_sgd.test(scr.test_X,params_meo_sc)
acc_test = me_sgd.evaluate(scr.test_y, y_pred_test)
10 print "Max-Ent Online Amazon Sentiment Accuracy train: %f
      test: %f"%(acc_train,acc_test)

```

## 1.5.2 Support Vector Machines

Support vector machines are also a discriminative approach, but they are not a probabilistic model at all. The basic idea is that, if the goal is to accurately predict outputs (according to some cost function), we should focus on that goal in the first place, rather than trying to estimate a probability distribution ( $P(Y|X)$  or  $P(X,Y)$ ), which is a more difficult problem. As Vapnik (1995) puts it, “do not solve an estimation problem of interest by solving a more general (harder) problem as an intermediate step.”

We next describe the *primal* problem associated with multi-class support vector machines (Crammer and Singer, 2002), which is of primary interest in natural language processing. There is a significant amount of literature about Kernel Methods (Schölkopf and Smola, 2002; Shawe-Taylor and Cristianini, 2004) mostly focused on the *dual* formulation. We will not discuss non-linear kernels or this dual formulation here.<sup>8</sup>

Consider  $\rho(y', y)$  as a non-negative cost function. For simplicity, we focus here on the 0/1-cost defined by Equation 1.26 (but keep in mind that other cost functions are possible). The *hinge loss*<sup>9</sup> is the function

$$\ell(\mathbf{w}; x, y) = \max_{y' \in \mathcal{Y}} \mathbf{w} \cdot \mathbf{f}(x, y') - \mathbf{w} \cdot \mathbf{f}(x, y) + \rho(y', y). \quad (1.34)$$

Note that the objective of Eq. 1.34 becomes zero when  $y' = y$ . Hence, we always have  $\ell(\mathbf{w}; x, y) \geq 0$ . Moreover, if  $\rho$  is the 0/1 cost, we have  $\ell(\mathbf{w}; x, y) = 0$  if and only if the weight vector is such that the model makes a correct prediction with a *margin* greater than 1: *i.e.*,  $\mathbf{w} \cdot \mathbf{f}(x, y) \geq \mathbf{w} \cdot \mathbf{f}(x, y') + 1$  for all  $y' \neq y$ . Otherwise, a positive loss is incurred.

<sup>8</sup>The main reason why we prefer to discuss the primal formulation with linear kernels is that the resulting algorithms run in linear time (or less), while known kernel-based methods are quadratic with respect to  $M$ . In large-scale problems (large  $M$ ) the former are thus more appealing.

<sup>9</sup>The hinge loss for the 0/1 cost is sometimes defined as  $\ell(\mathbf{w}; x, y) = \max\{0, \max_{y' \neq y} \mathbf{w} \cdot \mathbf{f}(x, y') - \mathbf{w} \cdot \mathbf{f}(x, y) + 1\}$ . Given our definition of  $\rho(\hat{y}, y)$ , note that the two definitions are equivalent.

Support vector machines (SVM) tackle the following optimization problem:

$$\hat{\boldsymbol{w}} = \arg \min_{\boldsymbol{w}} \sum_{m=1}^M \ell(\boldsymbol{w}; x^m, y^m) + \frac{\lambda}{2} \|\boldsymbol{w}\|^2, \quad (1.35)$$

where we also use the squared  $L_2$ -norm as the regularizer. For the 0/1-cost, the problem in Eq. 1.35 is equivalent to:

$$\arg \min_{\boldsymbol{w}, \xi} \sum_{m=1}^M \xi_m + \frac{\lambda}{2} \|\boldsymbol{w}\|^2 \quad (1.36)$$

$$\text{s.t. } \boldsymbol{w} \cdot \boldsymbol{f}(x^m, y^m) \geq \boldsymbol{w} \cdot \boldsymbol{f}(x^m, \tilde{y}^m) + 1 - \xi_m, \quad \forall m, \tilde{y}^m \in \mathcal{Y} \setminus \{y^m\} \quad (1.37)$$

Geometrically, we are trying to choose the linear classifier that yields the largest possible separation margin, while we allow some violations, penalizing the amount of slack via extra variables  $\xi_1, \dots, \xi_m$ .

Problem 1.35 does not have a closed form solution. Moreover, unlike maximum entropy models, here the objective function is non-differentiable, hence smooth optimization is not possible. However, it is still convex, which ensures that any local optimum is the global optimum. Despite not being differentiable, we can still define a *subgradient* of the objective function (which generalizes the concept of gradient), which enables us to apply subgradient-based methods. A stochastic subgradient algorithm for solving Eq. 1.35 is illustrated as Alg. 6. The similarity with maximum entropy models (Alg. 5) is striking: the only difference is that, instead of computing the feature vector expectation using the current model, we compute the feature vector associated with the cost-augmented prediction using the current model.

A variant of this algorithm was proposed by Shalev-Shwartz et al. (2007) under the name *Pegasos*, with excellent properties in large-scale settings. Other algorithms and software packages for training SVMs that have become popular are SVMLight (<http://svmlight.joachims.org>) and LIBSVM (<http://www.csie.ntu.edu.tw/~cjlin/libsvm/>), which allow non-linear kernels. These will generally be more suitable for smaller datasets, where high accuracy optimization can be obtained without much computational effort.

**Remark 1.3** Note the similarity between the stochastic (sub-)gradient algorithms (Algs. 5–6) and the online algorithms seen above (perceptron and MIRA).

**Exercise 1.7** Implement the SVM primal algorithm (Hint: look at the models implemented earlier, you should only need to change a few lines of code). Do this by creating a file `SVM.py` and implement class `SVM`. Then, repeat the MaxEnt exercise now using SVMs, for several values of  $\lambda$ :

```
import svm as svmc
2 svm = svmc.SVM()
```

---

**Algorithm 6** Stochastic Subgradient Descent for SVMs

---

- 1: **input:**  $\mathcal{D}$ ,  $\lambda$ , number of rounds  $T$ ,  
learning rate sequence  $(\eta_t)_{t=1,\dots,T}$
- 2: initialize  $w^1 = \mathbf{0}$
- 3: **for**  $t = 1$  **to**  $T$  **do**
- 4: choose  $m = m(t)$  randomly
- 5: take training pair  $(x^m, y^m)$  and compute the “cost-augmented prediction” under the current model:

$$\tilde{y} = \arg \max_{y' \in \mathcal{Y}} w^t \cdot f(x^m, y') - w^t \cdot f(x^m, y^m) + \rho(y', y)$$

- 6: update the model:

$$w^{t+1} \leftarrow (1 - \lambda \eta_t) w^t + \eta_t (f(x^m, y^m) - f(x^m, \tilde{y}))$$

- 7: **end for**
  - 8: **output:**  $\hat{w} \leftarrow w^{T+1}$
- 

```
4 svm.regularizer = 1.0 # This is lambda
  params_svm_sd = svm.train(sd.train_X, sd.train_y)
6 y_pred_train = svm.test(sd.train_X, params_svm_sd)
  acc_train = svm.evaluate(sd.train_y, y_pred_train)
8 y_pred_test = svm.test(sd.test_X, params_svm_sd)
  acc_test = svm.evaluate(sd.test_y, y_pred_test)
10 print "SVM Online Simple Dataset Accuracy train: %f test: %f" % (
    acc_train, acc_test)

12 fig, axis = sd.add_line(fig, axis, params_svm_sd, "SVM", "orange")

14 params_svm_sc = svm.train(scr.train_X, scr.train_y)
  y_pred_train = svm.test(scr.train_X, params_svm_sc)
16 acc_train = svm.evaluate(scr.train_y, y_pred_train)
  y_pred_test = svm.test(scr.test_X, params_svm_sc)
18 acc_test = svm.evaluate(scr.test_y, y_pred_test)
  print "SVM Online Amazon Sentiment Accuracy train: %f test: %f"
    %(acc_train, acc_test)
```

*Compare the results achieved and separating hiperplanes found.*

## 1.6 Comparison

Table 1.6 provides a high-level comparison among the different models discussed in this chapter.

	Naive Bayes	Perceptron	MIRA	MaxEnt	SVMs
Generative/Discriminative	G	D	D	D	D
Performance if true model not in the hypothesis class	Bad	Fair (may not converge)	Good	Good	Good
Performance if features overlap	Fair	Good	Good	Good	Good
Training	Closed Form	Easy	Easy	Fair	Fair
Hyperparameters to tune	1 (smoothing)	0	1	1	1

Table 1.1: Comparison among different models.

**Exercise 1.8** • *Using the simple data set run the different models varying some characteristics of the data, number of points, variance (hence separability), class balance. Use function XX which receives a dataset and plots all decisions boundaries and accuracies. What can you say about the methods when the amount of data increases? What about when the classes become too unbalanced.*

## 1.7 Final remarks

Some implementations of the discussed algorithms are available on the Web:

- SVMLight: <http://svmlight.joachims.org>
- LIBSVM: <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>
- Maximum Entropy: [http://homepages.inf.ed.ac.uk/lzhang10/maxent\\_toolkit.html](http://homepages.inf.ed.ac.uk/lzhang10/maxent_toolkit.html)
- MALLETT: <http://mallet.cs.umass.edu/>.