

## Lecture 7: Classification with Generative Algorithms

Adapted from Applied Machine Learning Lecture Notes of Volodymyr Kuleshov, Cornell Tech

**Instructor Tan Bui**

## Part 1: Generative Models

In this lecture, we are going to look at generative algorithms and their applications to classification.  
We will start by defining the concept of a generative *model*.

## Review: Components of A Supervised Machine Learning Problem

At a high level, a supervised machine learning problem has the following structure:

$$\underbrace{\text{Training Dataset} + \text{Attributes + Features}}_{\text{Learning Algorithm}} \rightarrow \text{Predictive Model}$$
$$\quad \quad \quad \underbrace{\text{Model Class + Objective + Optimizer}}$$

## Review: Probabilistic Models

A (parametric) probabilistic model with parameters  $\theta$  is a probability distribution

$$P_\theta(\mathbf{x}, y) : \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1].$$

This model can approximate the data distribution  $\mathbb{P}(\mathbf{x}, y)$ .

If we know  $P_\theta(\mathbf{x}, y)$ , we can compute predictions using the formula

$$P_\theta(y|\mathbf{x}) = \frac{P_\theta(\mathbf{x}, y)}{P_\theta(\mathbf{x})} = \frac{P_\theta(\mathbf{x}, y)}{\sum_{y \in \mathcal{Y}} P_\theta(\mathbf{x}, y)}.$$

Note that, for continuous random variable, we have

$$P_\theta(\mathbf{x}) = \int P_\theta(\mathbf{x}, y) dy$$

instead of the summation.

## Review: Maximum (joint) Likelihood Learning

In order to fit probabilistic models, we use the following objective:

$$\max_{\theta} \mathbb{E}_{\mathbf{x}, y \sim \mathbb{P}_{\text{data}}} \log P_{\theta}(\mathbf{x}, y).$$

This seeks to find a model that assigns high probability to the training data.

## Review: Maximum conditional Likelihood Learning

Alternatively, we may define a model of the conditional probability distribution:

$$P_{\theta}(y|\mathbf{x}) : \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1].$$

These are trained using conditional maximum likelihood:

$$\max_{\theta} \mathbb{E}_{\mathbf{x}, y \sim \mathbb{P}_{\text{data}}} \log P_{\theta}(y|\mathbf{x}).$$

This seeks to find a model that assigns high conditional probability to the target  $y$  for each  $\mathbf{x}$ .

Logistic regression is an example of this approach.

## Discriminative vs. Generative Models

These two types of models are also known as *generative* and *discriminative*.

$$\underbrace{P_{\theta}(\mathbf{x}, y) : \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]}_{\text{generative model}}$$
$$\underbrace{P_{\theta}(y|\mathbf{x}) : \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]}_{\text{discriminative model}}$$

- The models parametrize different kinds of probabilities
- They involve different training objectives and make different predictions
- Their uses are different (e.g., prediction, generation); more later!

$$\sim P_{\theta}(y|x) = \frac{P_{\theta}(x,y)}{P_{\theta}(x)}$$

can be used to predict (discriminate) what class  $x$  belongs to.

only allows us to calculate  $P_{\theta}(y|x)$

- once  $\theta$  is determined we can draw random samples  $(x, y)$  from the joint distribution  
 $P_{\theta}(x|y) = \frac{P_{\theta}(x,y)}{P_{\theta}(y)}$   
can be used to generate random samples given  $y$  ( $y = \text{setosa}$ )

## Classification Dataset: Iris Flowers

To demonstrate the two approaches, we are going to use the Iris flower dataset.

It's a classical dataset originally published by [R. A. Fisher](#) in 1936. Nowadays, it's widely used for demonstrating machine learning algorithms.

In [1]:

```
import numpy as np
import pandas as pd
import warnings
warnings.filterwarnings('ignore')
from sklearn import datasets

# Load the Iris dataset
iris = datasets.load_iris(as_frame=True)

# print part of the dataset
iris_X, iris_y = iris.data, iris.target
pd.concat([iris_X, iris_y], axis=1).head()
```

Out[1]:

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)	target
0	5.1	3.5	1.4	0.2	0
1	4.9	3.0	1.4	0.2	0
2	4.7	3.2	1.3	0.2	0
3	4.6	3.1	1.5	0.2	0
4	5.0	3.6	1.4	0.2	0

If we only consider the first two feature columns, we can visualize the dataset in 2D.

In [2]:

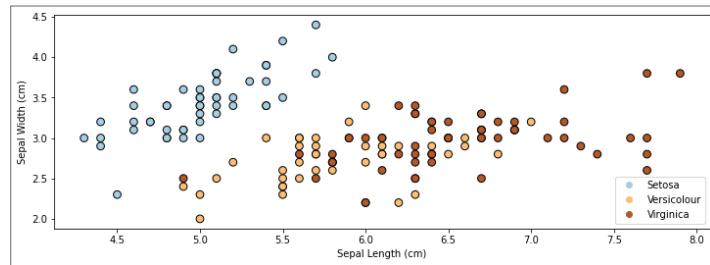
```
# https://scikit-learn.org/stable/auto_examples/neighbors/plot_classification.html
%matplotlib inline
from matplotlib import pyplot as plt
plt.rcParams['figure.figsize'] = [12, 4]

# create 2d version of dataset
X = iris_X.to_numpy()[:, :2]
x_min, x_max = X[:, 0].min() - .5, X[:, 0].max() + .5
y_min, y_max = X[:, 1].min() - .5, X[:, 1].max() + .5

# Plot also the training points
p1 = plt.scatter(X[:, 0], X[:, 1], c=iris_y, edgecolor='k', s=60, cmap=plt.cm.Paired)
plt.xlabel('Sepal Length (cm)')
plt.ylabel('Sepal Width (cm)')
plt.legend(handles=p1.legend_elements()[0], labels=['Setosa', 'Versicolour', 'Virginica'], loc='lower right')
```

Out[2]:

```
<matplotlib.legend.Legend at 0x7f2fa95f1710>
```



## Example: Discriminative Model

An example of a discriminative model is logistic or softmax regression.

- Discriminative models directly partition the feature space into regions associated with each class and separated by a decision boundary.
- Given features  $x$ , discriminative models directly map to predicted classes (e.g., via the function  $\sigma(\theta^\top x)$  for logistic regression).

In [3]:

```
# https://scikit-learn.org/stable/auto_examples/linear_model/plot_iris_logistic.html
from sklearn.linear_model import LogisticRegression
logreg = LogisticRegression(C=1e5, multi_class='multinomial')

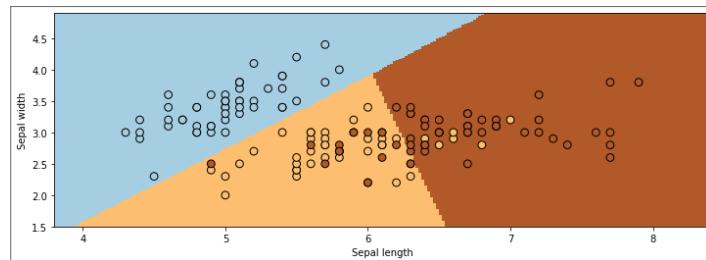
# Create an instance of Softmax and fit the data.
logreg.fit(X, iris_y)
xx, yy = np.meshgrid(np.arange(x_min, x_max, .02), np.arange(y_min, y_max, .02))
Z = logreg.predict(np.c_[xx.ravel(), yy.ravel()])

# Put the result into a color plot
Z = Z.reshape(xx.shape)
plt.pcolormesh(xx, yy, Z, cmap=plt.cm.Paired)

# Plot also the training points
plt.scatter(X[:, 0], X[:, 1], c=iris_y, edgecolors='k', s=60, cmap=plt.cm.Paired)
plt.xlabel('Sepal length')
plt.ylabel('Sepal width')
```

Out[3]:

Text(0, 0.5, 'Sepal width')



## Example: Generative Model

Generative modeling can be seen as taking a different approach:

1. In the Iris example, we first build a model of how each type of flower looks, i.e. we can learn the distribution

$$P(\mathbf{x}|y = k) \text{ for each class } k.$$

$$\frac{P(x, y)}{P(y)}$$

It defines a model of how each flower is generated, hence the name.

2. Given a new flower datapoint  $\mathbf{x}'$ , we can match it against each flower model  $P(\mathbf{x}|y)$  and find the type of flower  $y$  that looks most similar to it.

Mathematically, this corresponds to:

$$\begin{aligned} \arg \max_y \log P(y|\mathbf{x}) &= \arg \max_y \log \frac{P(\mathbf{x}|y)P(y)}{P(\mathbf{x})} \\ &= \arg \max_y \log P(\mathbf{x}|y)P(y), \end{aligned}$$

$$P(y|\mathbf{x}) = \frac{P(\mathbf{x}, y)}{P(\mathbf{x})}$$

where we have applied Bayes' rule in the first equality.

$$\left[ \begin{array}{l} P(y=1|\mathbf{x}') \\ P(y=2|\mathbf{x}') \\ \vdots \\ P(y=k|\mathbf{x}') \end{array} \right]$$

$\Leftarrow$  say it turns out that  
 $P(y=2|\mathbf{x}')$  is the largest  
 $\Rightarrow \mathbf{x}' \text{ belongs to the second class}$

## Generative vs. Discriminative Approaches

How do we know which approach is better?

- If we only care about prediction, we don't need a model of  $P(\mathbf{x})$ . We can solve precisely the problem we care about.
  - Discriminative models will often be more accurate.
- If we care about other tasks (generation, dealing with missing values, etc.) or if we know the true model is generative, we want to use the generative approach.

More on this later!

## Part 2: Gaussian Discriminant Analysis = a generative ML model

We are now going to continue our discussion of classification.

- We will see a new classification algorithm, Gaussian Discriminant Analysis.
- This will be our first example of generative machine learning model.

## Review: Classification

Consider a training dataset  $\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(n)}, y^{(n)})\}$ .

We distinguish between two types of supervised learning problems depending on the targets  $y^{(i)}$ .

1. **Regression:** The target variable  $y \in \mathcal{Y}$  is continuous:  $\mathcal{Y} \subseteq \mathbb{R}$ .

2. **Classification:** The target variable  $y$  is discrete and takes on one of  $K$  possible values:  $\mathcal{Y} = \{y_1, y_2, \dots, y_K\}$ . Each discrete value corresponds to a *class* that we want to predict.

## Review: Generative Models

There are two types of probabilistic models: *generative* and *discriminative*.

$$\underbrace{P_{\theta}(\mathbf{x}, y) : \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]}_{\text{generative model}} \quad \underbrace{P_{\theta}(y|\mathbf{x}) : \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]}_{\text{discriminative model}}$$

- They involve different training objectives and make different predictions
- Their uses are different (e.g., prediction, generation); more later!

## A brief introduction to Mixtures of Gaussians

A mixture of  $K$  Gaussians is a distribution  $P(\mathbf{x})$  of the form:

$$\phi_1 \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) + \phi_2 \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) + \dots + \phi_K \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K).$$

- Each  $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$  is a (multivariate) Gaussian distribution with mean  $\boldsymbol{\mu}_k$  and covariance  $\boldsymbol{\Sigma}_k$ .
- The  $\phi_k$  are weights, and the above sum is a weighted average of the  $K$  Gaussians.

gaussian distribution for  $\tilde{\mathbf{x}}$   
with mean  $\boldsymbol{\mu}_2$  and

covariance matrix

$$\boldsymbol{\Sigma}_2$$

when  $\sum_{i=1}^K \phi_i = 1$

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K)$$

We can easily visualize this in 1D:

In [4]:

```
def N(x,mu,sigma):
    return np.exp(-.5*(x-mu)**2/sigma**2)/np.sqrt(2*np.pi*sigma)

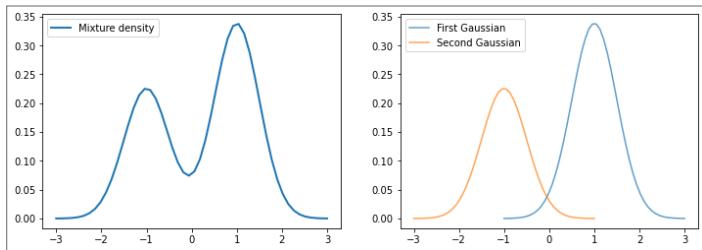
def mixture(x):
    return 0.6*N(x,mu=1,sigma=0.5) + 0.4*N(x,mu=-1,sigma=0.5)

xs, xs1, xs2 = np.linspace(-3,3), np.linspace(-1,3), np.linspace(-3,1)
plt.subplot('121')
plt.plot(xs, mixture(xs), label='Mixture density', linewidth=2)
plt.legend()

plt.subplot('122')
plt.plot(xs1, 0.6*N(xs1,mu=1,sigma=0.5), label='First Gaussian', alpha=0.7)
plt.plot(xs2, 0.4*N(xs2,mu=-1,sigma=0.5), label='Second Gaussian', alpha=0.7)
plt.legend()
```

Out[4]:

<matplotlib.legend.Legend at 0x7f2f9f770320>



## Gaussian Discriminant Model using mixture of Gaussians

We may use mixture of Gaussians to define a model  $P_\theta$ . This will be the basis of an algorithm called Gaussian Discriminant Analysis.

- The distribution over classes is Categorical, denoted  $\text{Categorical}(\phi_1, \phi_2, \dots, \phi_K)$ . We thus define  $P_\theta(y = k) = \phi_k$ .
- The conditional probability  $P_\theta(x | y = k)$  of the data under class  $k$  is a multivariate Gaussian  $\mathcal{N}(x; \mu_k, \Sigma_k)$  with mean and covariance  $\mu_k, \Sigma_k$ .

Thus, by the law of total probability (see the probability note),  $P_\theta(x)$  is a mixture of  $K$  Gaussians:

$$P_\theta(x) = \sum_{k=1}^K P_\theta(y = k) P_\theta(x | y = k) = \sum_{k=1}^K \phi_k \mathcal{N}(x; \mu_k, \Sigma_k)$$

assume



probability of  $x$  in  
class  $k$ .

probability of  $k$ th class

Intuitively, this model defines a story for how the data was generated. To obtain a data point,

- STEP 1:** we sample a class  $y \sim \text{Categorical}(\phi_1, \phi_2, \dots, \phi_K)$  with class proportions given by the  $\phi_k$ .
- STEP 2:** we sample an  $x$  from a Gaussian distribution  $\mathcal{N}(\mu_k, \Sigma_k)$  specific to that class.

Such a story can be constructed for most generative algorithms and helps understand them.

Note :  $P_\theta(y = k) P_\theta(x | y = k) = P_\theta(x, y = k)$

## Classification Dataset: Iris Flowers

To demonstrate this approach, we are going to use the Iris flower dataset.

It's a classical dataset originally published by [R. A. Fisher](#) in 1936. Nowadays, it's widely used for demonstrating machine learning algorithms.

In [5]:

```
import numpy as np
import pandas as pd
import warnings
warnings.filterwarnings('ignore')
from sklearn import datasets

# Load the Iris dataset
iris = datasets.load_iris(as_frame=True)

# print part of the dataset
iris_X, iris_y = iris.data, iris.target
pd.concat([iris_X, iris_y], axis=1).head()
```

Out[5]:

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)	target
0	5.1	3.5	1.4	0.2	0
1	4.9	3.0	1.4	0.2	0
2	4.7	3.2	1.3	0.2	0
3	4.6	3.1	1.5	0.2	0
4	5.0	3.6	1.4	0.2	0

If we only consider the first two feature columns, we can visualize the dataset in 2D.

In [6]:

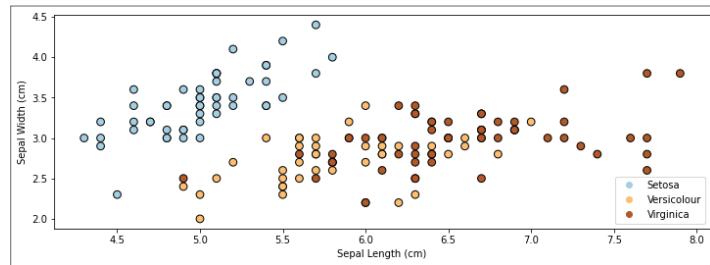
```
# https://scikit-learn.org/stable/auto_examples/neighbors/plot_classification.html
%matplotlib inline
from matplotlib import pyplot as plt
plt.rcParams['figure.figsize'] = [12, 4]

# create 2d version of dataset
X = iris_X.to_numpy()[:, :2]
x_min, x_max = X[:, 0].min() - .5, X[:, 0].max() + .5
y_min, y_max = X[:, 1].min() - .5, X[:, 1].max() + .5

# Plot also the training points
p1 = plt.scatter(X[:, 0], X[:, 1], c=iris_y, edgecolor='k', s=60, cmap=plt.cm.Paired)
plt.xlabel('Sepal Length (cm)')
plt.ylabel('Sepal Width (cm)')
plt.legend(handles=p1.legend_elements()[0], labels=['Setosa', 'Versicolour', 'Virginica'], loc='lower right')
```

Out[6]:

```
<matplotlib.legend.Legend at 0x7f2fa98eae48>
```



## Example: Iris Flower Classification

Let's see how this approach can be used in practice on the Iris dataset.

- We will "guess" a good set of parameters for a Gaussian Discriminant model
- We will sample from the model and compare to the true data

In [16]:

```
s = 100 # number of samples
K = 3 # number of classes
d = 2 # number of features

# guess the parameters
phi = 1./K * np.ones(K,)
mus = np.array(
[[5.0, 3.5],
[6.0, 2.5],
[6.5, 3.0]])
)
Sigmas = 0.05*np.tile(np.reshape(np.eye(2),(1,2,2)),(K,1,1))

# STEP 1: we sample flower class based on phi
ys = np.random.multinomial(n=1, pvals=phi, size=(s,)).argmax(axis=1)
xs = np.zeros([s,d])

# STEP 2: sample flowers from each class
for k in range(K):
nk = (ys==k).sum()
xs[ys==k,:] = np.random.multivariate_normal(mus[k], Sigmas[k], size=(nk,))

print(xs[:10])
```

```
[[5.2165882  3.21585548]
[6.33723498 2.37079596]
[4.70584828 3.66759626]
[6.01757646 2.87483419]
[5.63943149 2.31868567]
[6.12322516 2.75806014]
[6.3354125  2.81309308]
[6.05775836 2.5764983 ]
[6.02840033 2.80197148]
[6.11929254 2.43503439]]
```

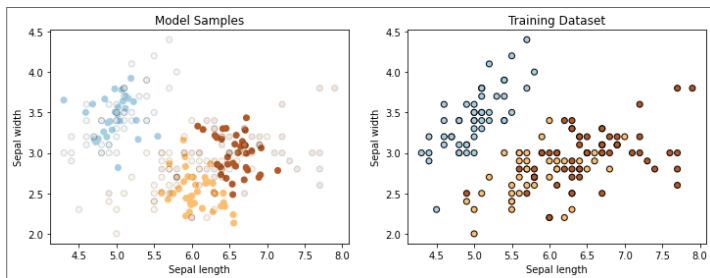
In [8]:

```
plt.subplot('121')
plt.title('Model Samples')
plt.scatter(xs[:,0], xs[:,1], c=ys, cmap=plt.cm.Paired)
plt.scatter(X[:, 0], X[:, 1], c=iris_y, edgecolors='k', cmap=plt.cm.Paired, alpha=0.15)
plt.xlabel('Sepal length')
plt.ylabel('Sepal width')

# Plot also the training points
plt.subplot('122')
plt.title('Training Dataset')
plt.scatter(X[:, 0], X[:, 1], c=iris_y, edgecolors='k', cmap=plt.cm.Paired, alpha=1)
plt.xlabel('Sepal length')
plt.ylabel('Sepal width')
```

Out[8]:

Text(0, 0.5, 'Sepal width')



- Our Gaussian Discirminant model generates data that looks not unlike the real data.
- Let's now see how we can learn parameters from data and use the model to make predictions.

## Part 3: Gaussian Discriminant Analysis: Learning

In the previous example, we show the two-step process to generate new flowers provided that we know the parameters ,  $\{\phi_i\}_{i=1}^K, \{\boldsymbol{\mu}_i\}_{i=1}^K, \{\boldsymbol{\Sigma}_i\}_{i=1}^K$ . We continue our discussion of Gaussian Discriminant analysis, and look at:

- How to learn these parameters of the mixture model
- How to use the model to make predictions

## Review: Gaussian Discriminant Model

We may define a model  $P_{\theta}$  as follows. This will be the basis of an algorithm called Gaussian Discriminant Analysis.

- The distribution over classes is **Categorical**, denoted  $\text{Categorical}(\phi_1, \phi_2, \dots, \phi_K)$ . Thus,  $P_{\theta}(y = k) = \phi_k$ .
- The conditional probability  $P(\mathbf{x} | y = k)$  of the data under class  $k$  is a **multivariate Gaussian**  $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$  with mean and covariance  $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$ .

Thus,  $P_{\theta}(\mathbf{x})$  is a mixture of  $K$  Gaussians:

$$P_{\theta}(\mathbf{x}) = \sum_{k=1}^K P_{\theta}(y = k) P_{\theta}(\mathbf{x}|y = k) = \sum_{k=1}^K \phi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k).$$

The **joint** distribution of  $\mathbf{x}$  and  $y = i$  is thus

$$P_{\theta}(\mathbf{x}, y = i) = P_{\theta}(\mathbf{x}|y = i) P_{\theta}(y = i) = \phi_i \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i),$$

and this is exactly what we need for the joint likelihood next.

## Review: Maximum Likelihood Learning

In order to fit probabilistic models, we use maximum joint likelihood that we have learned:

$$\max_{\theta} \mathbb{E}_{\mathbf{x}, y \sim \mathbb{P}_{\text{data}}} \log P_{\theta}(\mathbf{x}, y).$$

This seeks to find a model that assigns high probability to the training data.

Let's use maximum likelihood to fit the Gaussian Discriminant model. In this case, the model parameters  $\theta$  are the union of the parameters of each sub-model:

$$\theta = (\mu_1, \Sigma_1, \phi_1, \dots, \mu_K, \Sigma_K, \phi_K).$$

$$\begin{aligned} & \Leftrightarrow P_{\theta}(\mathbf{x}, y = i) \\ &= \phi_i \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \end{aligned}$$

As discuss above, the components of the joint likelihood  $P_{\theta}(\mathbf{x}, y = i)$  are as follows:

$$P_{\theta}(y = i) = \frac{\phi_i}{\sum_{k=1}^K \phi_k}$$

$$P_{\theta}(\mathbf{x}|y = i) = \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \exp\left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^\top \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i)\right).$$

Note that we have used  $\frac{\phi_i}{\sum_{k=1}^K \phi_k}$  instead of  $\phi_i$  to avoid constrained optimization. Indeed, using  $P_{\theta}(y = i) = \phi_i$  requires the constraint  $\sum_{k=1}^K \phi_k = 1$ . While it is possible to deal with constrained optimization (such as the method of Lagrangian multipliers), it is typically harder than unconstrained optimization.

$$\Rightarrow P_{\theta}(\mathbf{x}, y = i) \neq \phi_i;$$

$$= \frac{\phi_i}{\sum_{k=1}^K \phi_k}$$

## Optimizing the Log Likelihood

In the following, we use  $P_\beta(\cdot) = P(\cdot|\beta)$ , where  $\beta$  is a set of parameters such as  $(\mu_j, \Sigma_j)$  or  $\phi = \{\phi_1, \dots, \phi_K\}$  or the whole set of training parameters  $\theta$ . Given a I.I.D. dataset  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)}) \mid i = 1, 2, \dots, n\}$ , we want to optimize the log-likelihood  $\ell(\theta)$ :

$$\begin{aligned}\ell(\theta) &= \log \prod_{i=1}^n P_\theta(\mathbf{x}^{(i)}, y^{(i)}) = \sum_{i=1}^n \log P_\theta(\mathbf{x}^{(i)}, y^{(i)}) = \sum_{i=1}^n \log P_\theta(\mathbf{x}^{(i)}|y^{(i)}) + \sum_{i=1}^n \log P_\theta(y^{(i)}) \\ &= \sum_{j=1}^K \underbrace{\sum_{i:y^{(i)}=j} \log P(\mathbf{x}^{(i)}|y^{(i)}; \mu_j, \Sigma_j)}_{\text{all the terms that involve } \mu_k, \Sigma_k} + \underbrace{\sum_{i=1}^n \log P(y^{(i)}|\phi)}_{\text{all the terms that involve } \phi} \stackrel{\text{"}}{=} \ell_\phi(y^{(i)}) \\ &\quad \text{Note: } \ell_\phi(y^{(i)}) \text{ does not depend on } \mu_k, \Sigma_k \\ &\quad \phi = (\phi_1, \dots, \phi_K)^\top\end{aligned}$$

Notice that each set of parameters  $(\mu_j, \Sigma_j)$  is found in only the first term  $\sum_{i:y^{(i)}=j} \log P(\mathbf{x}^{(i)}|y^{(i)}; \mu_j, \Sigma_j)$  of the summation over the  $K$  classes and the  $\phi$  are in the second term  $\sum_{i=1}^n \log P(y^{(i)}|\phi)$ .

Note that optimization with separable objective function such as  $\max_{z_1, \dots, z_K} \sum_{j=1}^K f(z_j)$  is equivalent with  $K$  independent optimization  $\max_{z_j} f(z_j)$ . Since each  $(\mu_j, \Sigma_j)$  for  $j = 1, 2, \dots, K$  is found in one term, optimization over  $(\mu_j, \Sigma_j)$  can be carried out independently:

$$\max_{\mu_j, \Sigma_j} \sum_{i:y^{(i)}=j} \log P(\mathbf{x}^{(i)} | y^{(i)}; \mu_j, \Sigma_j).$$

Similarly, optimizing for  $\phi = (\phi_1, \phi_2, \dots, \phi_K)$  only involves the second term:

$$\max_{\phi} \sum_{i=1}^n \log P(y^{(i)} | \phi)$$

## Optimizing the Class Probabilities

The above observations greatly simplify the optimization of the model. Let's first consider the optimization over  $\phi = (\phi_1, \phi_2, \dots, \phi_K)$ . From the previous analysis, our objective  $J(\phi)$  equals

$$\begin{aligned}
 J(\phi) &= \sum_{i=1}^n \log P_{\phi}(y^{(i)} | \phi) \\
 &= \sum_{i=1}^n \log \phi_{y^{(i)}} - n \cdot \log \sum_{k=1}^K \phi_k \\
 &= \sum_{k=1}^K \sum_{i:y^{(i)}=k} \log \phi_k - n \cdot \log \sum_{k=1}^K \phi_k
 \end{aligned}$$

$$P(y^{(i)} | \vec{\phi}) = \frac{\phi_{y^{(i)}}}{\sum_{k=1}^K \phi_k}$$

$$\sum_{i=1}^n \log P(y^{(i)} | \vec{\phi}) = \sum_{i=1}^n \log \frac{\phi_{y^{(i)}}}{\sum_{k=1}^K \phi_k}$$

$$y^{(i)} \in \{1, \dots, K\}$$

Taking the partial derivatives

$$\begin{aligned}
\nabla_{\phi_l} J(\boldsymbol{\phi}) &= \nabla_{\phi_l} \left( \sum_{k=1}^K \sum_{i:y^{(i)}=k} \log \phi_k - n \cdot \log \sum_{k=1}^K \phi_k \right) \\
&= \sum_{i:y^{(i)}=l} \frac{1}{\phi_l} - n \frac{1}{\sum_{k=1}^K \phi_k} \\
&= \frac{n_l}{\phi_l} - \frac{n}{\sum_{k=1}^K \phi_k}
\end{aligned}$$

Setting the derivative to zero (we switch dummy variables), we obtain

$$\frac{\phi_k}{\sum_l \phi_l} = \frac{n_k}{n}$$

for each  $k$ , where  $n_k = |\{i : y^{(i)} = k\}|$  is the number of training targets with class  $k$ .

Thus, the optimal  $\phi_k$  is just the proportion of data points with class  $k$  in the training set!

## Optimizing Conditional Probabilities

Similarly, we can maximize the likelihood

$$\max_{\mu_k, \Sigma_k} \sum_{i:y^{(i)}=k} \log P(\mathbf{x}^{(i)} | y^{(i)}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \max_{\mu_k, \Sigma_k} \sum_{i:y^{(i)}=k} \log \mathcal{N}(\mathbf{x}^{(i)} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

over the Gaussian parameters.

Computing the derivative and setting it to zero, we obtain closed form solutions:

$$\boldsymbol{\mu}_k = \frac{\sum_{i:y^{(i)}=k} \mathbf{x}^{(i)}}{n_k}$$

$$\boldsymbol{\Sigma}_k = \frac{\sum_{i:y^{(i)}=k} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_k)(\mathbf{x}^{(i)} - \boldsymbol{\mu}_k)^\top}{n_k}$$

→ empirical mean  
→ empirical covariance

These are just the empirical means and covariances of each class.

matrix

## Querying the Model

How do we ask the model for predictions? As discussed earlier, we can apply Bayes' rule:

$$\arg \max_y P_{\theta}(y|\mathbf{x}) = \arg \max_y P_{\theta}(\mathbf{x}|y)P(y).$$

Thus, we can estimate the probability of  $x$  under each  $P_{\theta}(\mathbf{x}|y = k)P(y = k)$  and choose the class that explains the data best.

## Classification Dataset: Iris Flowers

To demonstrate this approach, we are going to use the Iris flower dataset.

It's a classical dataset originally published by [R. A. Fisher](#) in 1936. Nowadays, it's widely used for demonstrating machine learning algorithms.

In [9]:

```
import numpy as np
import pandas as pd
import warnings
warnings.filterwarnings('ignore')
from sklearn import datasets

# Load the Iris dataset
iris = datasets.load_iris(as_frame=True)

# print part of the dataset
iris_X, iris_y = iris.data, iris.target
pd.concat([iris_X, iris_y], axis=1).head()
```

Out[9]:

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)	target
0	5.1	3.5	1.4	0.2	0
1	4.9	3.0	1.4	0.2	0
2	4.7	3.2	1.3	0.2	0
3	4.6	3.1	1.5	0.2	0
4	5.0	3.6	1.4	0.2	0

If we only consider the first two feature columns, we can visualize the dataset in 2D.

In [10]:

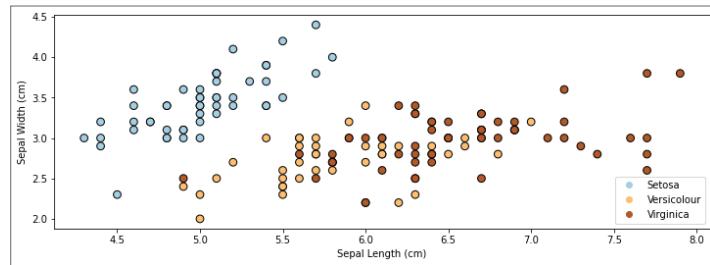
```
# https://scikit-learn.org/stable/auto_examples/neighbors/plot_classification.html
%matplotlib inline
from matplotlib import pyplot as plt
plt.rcParams['figure.figsize'] = [12, 4]

# create 2d version of dataset
X = iris_X.to_numpy()[:, :2]
x_min, x_max = X[:, 0].min() - .5, X[:, 0].max() + .5
y_min, y_max = X[:, 1].min() - .5, X[:, 1].max() + .5

# Plot also the training points
p1 = plt.scatter(X[:, 0], X[:, 1], c=iris_y, edgecolor='k', s=60, cmap=plt.cm.Paired)
plt.xlabel('Sepal Length (cm)')
plt.ylabel('Sepal Width (cm)')
plt.legend(handles=p1.legend_elements()[0], labels=['Setosa', 'Versicolour', 'Virginica'], loc='lower right')
```

Out[10]:

```
<matplotlib.legend.Legend at 0x7f2f9f6f2358>
```



## Example: Iris Flower Classification

Let's see how this approach can be used in practice on the Iris dataset.

- We will learn a good set of parameters for a Gaussian Discriminant model
- We will compare the outputs to the true predictions.

Let's first start by computing the true parameters on our dataset.

In [11]:

```
# we can implement these formulas over the Iris dataset
d = 2 # number of features in our toy dataset
K = 3 # number of classes
n = X.shape[0] # size of the dataset

# these are the shapes of the parameters
mus = np.zeros([K,d])
Sigmas = np.zeros([K,d,d])
phis = np.zeros([K])

# we now compute the parameters from the solution of the above optimization
for k in range(3):
    X_k = X[iris_y == k]
    mus[k] = np.mean(X_k, axis=0)
    Sigmas[k] = np.cov(X_k.T)
    phis[k] = X_k.shape[0] / float(n)

# print out the means
print(mus)
```

```
[[5.006 3.428]
 [5.936 2.77 ]
 [6.588 2.974]]
```

We can compute predictions using Bayes' rule.

In [12]:

```
# we can implement this in numpy
def gda_predictions(x, mus, Sigmas, phis):
    """This returns class assignments and p(y|x) under the GDA model.

    We compute \arg\max_y p(y|x) as \arg\max_y p(x|y)p(y)
    """

    # adjust shapes
    n, d = x.shape
    x = np.reshape(x, (1, n, d, 1))
    mus = np.reshape(mus, (K, 1, d, 1))
    Sigmas = np.reshape(Sigmas, (K, 1, d, d))

    # compute probabilities
    py = np.tile(phis.reshape((K,1)), (1,n)).reshape([K,n,1,1])
    pxy = (
        np.sqrt(np.abs((2*np.pi)**d*np.linalg.det(Sigmas))).reshape((K,1,1,1))
        * -.5*np.exp(
            np.matmul(np.matmul((x-mus).transpose([0,1,3,2]), np.linalg.inv(Sigmas)), x-mus)
        )
    )
    pyx = pxy * py
    return pyx.argmax(axis=0).flatten(), pyx.reshape([K,n])

idx, pyx = gda_predictions(X, mus, Sigmas, phis)
print(idx)
```

```
[0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 2 2 2 2 1 2 1 2 1 2 1 1 1 1 1 1 1 1 2 1 1 1 1 1 1 1 1 2 1
 2 2 2 2 2 1 1 1 1 1 1 1 2 2 2 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 1 2 2 2 2 1 2 2 2 2 2 2
 2 2 1 1 2 2 2 2 1 2 1 2 2 2 1 1 2 2 2 2 1 1 2 2 2 1 2 2 2 1 2 2 2 1 2 2 2 2 2 2 2
 2 1]
```

We visualize the decision boundaries like we did earlier.

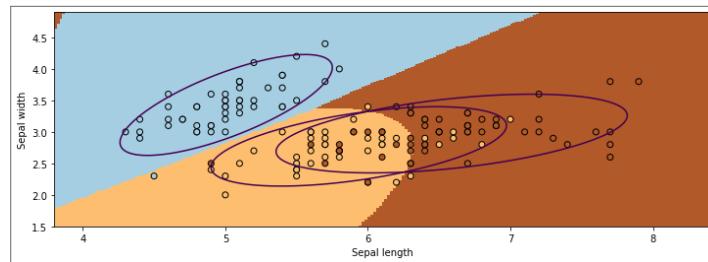
In [13]:

```
from matplotlib.colors import LogNorm
xx, yy = np.meshgrid(np.arange(x_min, x_max, .02), np.arange(y_min, y_max, .02))
Z, pyx = gda_predictions(np.c_[xx.ravel(), yy.ravel()], mus, Sigmas, phis)
logpy = np.log(-1./3*pyx)

# Put the result into a color plot
Z = Z.reshape(xx.shape)
contours = np.zeros([K, xx.shape[0], xx.shape[1]])
for k in range(K):
    contours[k] = logpy[k].reshape(xx.shape)
plt.pcolormesh(xx, yy, Z, cmap=plt.cm.Paired)
for k in range(K):
    plt.contour(xx, yy, contours[k], levels=np.logspace(0, 1, 1))

# Plot also the training points
plt.scatter(X[:, 0], X[:, 1], c=iris_y, edgecolors='k', cmap=plt.cm.Paired)
plt.xlabel('Sepal length')
plt.ylabel('Sepal width')

plt.show()
```



## Algorithm: Gaussian Discriminant Analysis

- **Type:** Supervised learning (multi-class classification)
- **Model family:** Mixtures of Gaussians.
- **Objective function:** Log-likelihood.
- **Optimizer:** Closed form solution.

## Special Cases of GDA

Many important generative algorithms are special cases of Gaussian Discriminative Analysis

- Linear discriminant analysis (LDA): all the covariance matrices  $\Sigma_k$  take the same value.
- Gaussian Naive Bayes: all the covariance matrices  $\Sigma_k$  are diagonal.
- Quadratic discriminant analysis (QDA): another term for GDA.

→ decision boundary is linear

## Generative vs. Discriminative Approaches

Discriminative classifiers Goal: Want to predict class probability given  $x$ , that is, we want to compute  $P_{\theta}(y|x)$ . We accomplish this through two steps

- Assume a model form for  $P_{\theta}(y|x)$
- Train  $\theta^*$  using maximum likelihood (or Bayesian approach) using training data.

*conditional*

Pros of discriminative models:

- Often more accurate because they make fewer modeling assumptions.

Generative classifier Goal: Want to predict class probability given  $x$ , that is, we want to compute  $P_{\theta}(y|x)$ . We accomplish that using Bayesian approach:

$$P_{\theta}(y|x) = \frac{P_{\theta}(\mathbf{x}|y)P_{\theta}(y)}{P_{\theta}(\mathbf{x})}$$

via the following steps

- Postulate a model for class probability:  $P_{\theta}(y)$
- Postulate a model for the conditional probability of feature  $x$  within a class:  $P_{\theta}(\mathbf{x}|y)$
- Train  $\theta^*$  based on the joint likelihood  $P_{\theta}(\mathbf{x}, y) = P_{\theta}(\mathbf{x}|y)P_{\theta}(y)$

Pros of generative models:

1. Can do more than just prediction: generation, fill-in missing features, etc.
2. Can include extra prior knowledge; if prior knowledge is correct, model will be more accurate.
3. Often have closed-form solutions, hence are faster to train.