

Tutorial Linear Classification & Kernels

Suggested reading for linear classification:

Murphy: 4.2 (examples for using gaussians as class conditional densities) (optional), 8.1, 8.2, 8.3, 8.6 until 8.6.1 (incl.) and Barber 17.4 (part on the perceptron).
as always very elegant and concise, contains a lot of links and further reading (e.g. on bayesian logistic regression etc.)

OR

Barber: 17.4, 17.6
very concise, the whole chapter 17 covers lin. regression, lin. Classification, SVMs, Kernels and optimization techniques in a very economical easy to understand and short way.
probably the reading of choice for all of these topics if time is scarce and/or passing the course is the main interest.

OR

Bishop: 4.1.1, 4.1.2, 4.1.7, 4.2, 4.3.2, 4.3.4
well rounded presentation. as always didactically nice and complete.

Suggested Reading on Kernels:

Murphy ch 14

OR

Bishop ch 6

OR

Barber, chapter 17

Part I: Problems w.r.t. Kernels

1 The Gaussian kernel

Problem 1. One of the nice things about kernels is that new kernels can be constructed out of already given ones. Use the five kernel construction rules from the lecture to prove that the function

$$\exp\left(-\frac{|\boldsymbol{x} - \boldsymbol{y}|^2}{2\sigma^2}\right)$$

is a kernel.

(Hint: Use the Taylor expansion of the exponential function to prove that $\exp(K_1(\boldsymbol{x}, \boldsymbol{y}))$ is a kernel if $K_1(\boldsymbol{x}, \boldsymbol{y})$ is a kernel.)

2nd hint: It might help to apply the rule $K_3(\phi(\boldsymbol{x}), \phi(\boldsymbol{y}))$ with the linear kernel $K_3(\boldsymbol{x}, \boldsymbol{y}) = \boldsymbol{x}^T \boldsymbol{y}$ and consider a feature map $\phi(\boldsymbol{z})$ with only one feature.

Let K_1 and K_2 be kernels on $\mathcal{X} \subseteq \mathbb{R}^n$, then the following functions are kernels:

1. $K(\mathbf{x}, \mathbf{y}) = K_1(\mathbf{x}, \mathbf{y}) + K_2(\mathbf{x}, \mathbf{y})$,

2. $K(\mathbf{x}, \mathbf{y}) = a K_1(\mathbf{x}, \mathbf{y})$ for $a > 0$,

3. $K(\mathbf{x}, \mathbf{y}) = K_1(\mathbf{x}, \mathbf{y})K_2(\mathbf{x}, \mathbf{y})$,

4. $K(\mathbf{x}, \mathbf{y}) = K_3(\boldsymbol{\phi}(\mathbf{x}), \boldsymbol{\phi}(\mathbf{y}))$ for K_3 kernel on \mathbb{R}^m and $\boldsymbol{\phi}: \mathcal{X} \rightarrow \mathbb{R}^m$,

5. $K(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{B} \mathbf{y}$ for \mathbf{B} symmetric and positive semi-definite $n \times n$ matrix.

Let K_1 and K_2 be kernels on $\mathcal{X} \subseteq \mathbb{R}^n$, then the following functions are kernels:

1. $K(x, y) = K_1(x, y) + K_2(x, y)$,
 2. $K(x, y) = a K_1(x, y)$ for $a > 0$,
 3. $K(x, y) = K_1(x, y)K_2(x, y)$,
 4. $K(x, y) = K_3(\phi(x), \phi(y))$ for K_3 kernel on \mathbb{R}^m and $\phi: \mathcal{X} \rightarrow \mathbb{R}^m$,
 5. $K(x, y) = x^T B y$ for B symmetric and positive semi-definite $n \times n$ matrix.
-

$$\exp(K_1(x, y)) = 1 + \sum_{n=1}^{\infty} \frac{1}{n!} K_1(x, y)^n$$

$K_1(x, y)^n$ is a kernel R3

$(1/n!)K_1(x, y)^n$ is a kernel R2

$\sum_{n=1}^{\infty} 1/(n!)K_1(x, y)^n$ is a kernel R1

1 is a kernel : R5: $x^T y$ is kernel ($B=\text{Id.}$) (with $x, y \in \mathbb{R}^m$)

R4: $\Phi(x) = \Phi(y) = (1/\sqrt{m}, \dots, 1/\sqrt{m})$

or:

R5: $x^T y$ is kernel ($B=\text{Id.}$) (with $x, y \in \mathbb{R}^{m=1}$)

R4: $\Phi(x) = \Phi(y) = 1$

Let K_1 and K_2 be kernels on $\mathcal{X} \subseteq \mathbb{R}^n$, then the following functions are kernels:

1. $K(\mathbf{x}, \mathbf{y}) = K_1(\mathbf{x}, \mathbf{y}) + K_2(\mathbf{x}, \mathbf{y})$,
 2. $K(\mathbf{x}, \mathbf{y}) = a K_1(\mathbf{x}, \mathbf{y})$ for $a > 0$,
 3. $K(\mathbf{x}, \mathbf{y}) = K_1(\mathbf{x}, \mathbf{y})K_2(\mathbf{x}, \mathbf{y})$,
 4. $K(\mathbf{x}, \mathbf{y}) = K_3(\boldsymbol{\phi}(\mathbf{x}), \boldsymbol{\phi}(\mathbf{y}))$ for K_3 kernel on \mathbb{R}^m and $\boldsymbol{\phi}: \mathcal{X} \rightarrow \mathbb{R}^m$,
 5. $K(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{B} \mathbf{y}$ for \mathbf{B} symmetric and positive semi-definite $n \times n$ matrix.
-

$$\begin{aligned} \exp\left(-\frac{|\mathbf{x} - \mathbf{y}|^2}{2\sigma^2}\right) &= \exp\left(-\frac{(\mathbf{x} - \mathbf{y})^T (\mathbf{x} - \mathbf{y})}{2\sigma^2}\right) = \exp\left(-\frac{\mathbf{x}^T \mathbf{x} - 2\mathbf{x}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}}{2\sigma^2}\right) \\ &= \exp\left(-\frac{\mathbf{x}^T \mathbf{x}}{2\sigma^2}\right) \exp\left(-\frac{\mathbf{y}^T \mathbf{y}}{2\sigma^2}\right) \exp\left(\frac{\mathbf{x}^T \mathbf{y}}{\sigma^2}\right). \end{aligned}$$

R3

R5,R2

R5: $\mathbf{x}^T \mathbf{y}$ is kernel ($\mathbf{B} = \text{Id.}$), with $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{m=1}$

R4: $\Phi(\mathbf{x}) = \exp\left(\frac{-\mathbf{x}^T \mathbf{x}}{2\sigma}\right)$, $\Phi(\mathbf{y}) = \exp\left(\frac{-\mathbf{y}^T \mathbf{y}}{2\sigma}\right)$

slightly alternate approach to problem 1:

Problem 1a: Verify that $k(X_i, X_j) = f(X_i)k_1(X_i, X_j)f(X_j)$, where k_1 is a valid kernel, is also a kernel.

Problem 1a: Verify that $k(X_i, X_j) = f(X_i)k_1(X_i, X_j)f(X_j)$, where k_1 is a valid kernel, is also a kernel.

A kernel is defined as $k(X_i, X_j) = \phi(X_i)^T \phi(X_j)$. I.e. it is the scalar product between the feature vectors for two input vectors. If k_1 is a valid kernel then it follows that

$$ck_1(X_i, X_j) = u(X_i)^T u(X_j)$$

where

$$u(X_n) = c^{1/2} \phi(X_n)$$

and so $ck_1(X_i, X_j)$ can be expressed as the scalar product of feature vectors. This is how we show the validity of rule 2. We can follow a similar path for showing the above function is also a kernel.

We can write

$$f(X_i)k_1(X_i, X_j)f(X_j) = v(X_i)^T v(X_j)$$

where we have defined $v(X_n) = f(X_n)\phi(X_n)$.

Again, we see that $f(X_i)k_1(X_i, X_j)f(X_j)$ can be expressed as the scalar product of feature vectors, and hence is a valid kernel.

Problem 1b: Verify that $k(X_i, X_j) = q(k_1(X_i, X_j))$ (where q is a polynomial and k_1 is a valid kernel) and $k(X_i, X_j) = \exp(k_1(X_i, X_j))$ are valid rules for constructing a valid kernel.

Problem 1b: Verify that $k(X_i, X_j) = q(k_1(X_i, X_j))$ (where q is a polynomial and k_1 is a valid kernel) and $k(X_i, X_j) = \exp(k_1(X_i, X_j))$ are valid rules for constructing a valid kernel.

We can show that $k(X_i, X_j) = q(k_1(X_i, X_j))$ follows from $k(x, y) = c(k_1(X_i, X_j))$ where $c > 0$ is a constant (Rule 2), $k(X_i, X_j) = (k_1(X_i, X_j)k_2(X_i, X_j))$ where k_2 is also a valid kernel (Rule 3), and $k(X_i, X_j) = k_1(X_i, X_j) + k_2(X_i, X_j)$ (Rule 1).

For $k(X_i, X_j) = \exp(k_1(X_i, X_j))$ we re-write the exponential as a power series yielding:

$$\begin{aligned} k(X_i, X_j) &= \exp(k_1(X_i, X_j)) \\ &= \sum_{m=0}^{\infty} \frac{(k_1(X_i, X_j))^m}{m!} \end{aligned}$$

Since this is a polynomial in $k_1(X_i, X_j)$ we can use our proof from above to show that this is also a valid kernel.

Problem 1c: One commonly used kernel is the Gaussian kernel, i.e.

$$k(X_i, X_j) = \exp\left(\frac{-\|X_i - X_j\|^2}{2\sigma^2}\right),$$

but remember that in this context it is not interpreted as a probability density, that's why there isn't a normalization coefficient. We can see that this is a valid kernel by expanding the square

$$\|X_i - X_j\|^2 = X_i^T X_i + X_j^T X_j + 2X_i^T X_j$$

which gives us

$$k(X_i, X_j) = \exp\left(\frac{-X_i^T X_i}{2\sigma^2}\right) \exp\left(\frac{X_i^T X_j}{\sigma^2}\right) \exp\left(\frac{-X_j^T X_j}{2\sigma^2}\right)$$

Which we can derive since we know $k(X_i, X_j) = f(X_i)k_1(X_i, X_j)f(X_j)$ is a valid kernel, and we know that $k(X_i, X_j) = \exp(k_1(X_i, X_j))$ is a valid kernel.

So, by using the expanded version of the kernel from above, and expanding the middle factor as a power series, show that the Gaussian kernel equation we showed at the very top of this exercise can be expressed as the inner product of an infinite-dimensional feature vector.

Since we know that the exponential kernel, $k(X_i, X_j) = \exp(k_1(X_i, X_j))$, can be written as the infinite sum of terms, each of which can itself be written as an inner product of feature vectors (since we know that $k(X_i, X_j) = q(k_1(X_i, X_j))$ where q is a polynomial, is a valid kernel). Thus, by concatenating the feature vectors of the individual terms in that sum, we can write this as an inner product of infinite dimension vectors. More formally,

$$\begin{aligned} \exp(X_i^T X_j / \sigma^2) &= \sum_{m=0}^{\infty} \phi_m(X_i)^T \phi_m(X_j) \\ &= \psi(X_i)^T \psi(X_j) \end{aligned}$$

where $\psi(X_n)^T = [\phi_0(X_n)^T, \phi_1(X_n)^T, \dots]$. Hence, we can write the original form for the Gaussian kernel as

$$k(X_i, X_j) = \rho(X_i)^T \rho(X_j)$$

where

$$\rho(X_n) = \exp\left(\frac{X_n^T X_n}{\sigma^2}\right) \psi(X_n).$$

end of slightly alternate approach to problem 1.

Problem 2: Find an infinite-dimensional feature space $\phi(\mathbf{x})$ corresponding to the Gaussian kernel, i.e. determine $\phi(\mathbf{x})$ so that

$$\phi(\mathbf{x})^T \phi(\mathbf{y}) = \exp\left(-\frac{|\mathbf{x} - \mathbf{y}|^2}{2\sigma^2}\right) .$$

(Hint: The multinomial formula turns a power of a sum into a weighted sum of products,

$$\left(\sum_{t=1}^m x_t\right)^n = \sum_{k_1+k_2+\dots+k_m=n} \binom{n}{k_1, k_2, \dots, k_m} \prod_{t=1}^m x_t^{k_t} ,$$

with $\binom{n}{k_1, k_2, \dots, k_m} = \frac{n!}{k_1! k_2! \dots k_m!}$.)

Let $\mathbf{x}, \mathbf{y} \in \mathbb{R}^m$ and without loss of generality let $\sigma = 1$. We expand the Gaussian kernel into an infinite sum using the Taylor expansion of the exponential function and use the multinomial formula to separate the scalar product,

$$\begin{aligned}
\exp\left(-\frac{|\mathbf{x} - \mathbf{y}|^2}{2}\right) &= \exp\left(-\frac{\mathbf{x}^2}{2}\right) \exp\left(-\frac{\mathbf{y}^2}{2}\right) \exp(-\mathbf{x}^T \mathbf{y}) \\
&= \exp\left(-\frac{\mathbf{x}^2}{2}\right) \exp\left(-\frac{\mathbf{y}^2}{2}\right) \sum_{n=0}^{\infty} \frac{(\mathbf{x}^T \mathbf{y})^n}{n!} \\
&= \exp\left(-\frac{\mathbf{x}^2}{2}\right) \exp\left(-\frac{\mathbf{y}^2}{2}\right) \sum_{n=0}^{\infty} \frac{1}{n!} \left(\sum_{t=1}^m x_t y_t\right)^n \\
&= \exp\left(-\frac{\mathbf{x}^2}{2}\right) \exp\left(-\frac{\mathbf{y}^2}{2}\right) \sum_{n=0}^{\infty} \frac{1}{n!} \sum_{k_1+k_2+\dots+k_m=n} \binom{n}{k_1, k_2, \dots, k_m} \prod_{t=1}^m (x_t y_t)^{k_t} \\
&= \sum_{n=0}^{\infty} \sum_{k_1+k_2+\dots+k_m=n} \frac{\exp\left(-\frac{\mathbf{x}^2}{2}\right)}{\sqrt{n!}} \sqrt{\binom{n}{k_1, k_2, \dots, k_m}} \left(\prod_{t=1}^m x_t^{k_t}\right) \\
&\quad \frac{\exp\left(-\frac{\mathbf{y}^2}{2}\right)}{\sqrt{n!}} \sqrt{\binom{n}{k_1, k_2, \dots, k_m}} \left(\prod_{t=1}^m y_t^{k_t}\right).
\end{aligned}$$

Thus the feature space of the Gaussian kernel is

$$\phi(\mathbf{x}) = \left[\frac{\exp\left(-\frac{\mathbf{x}^2}{2}\right)}{\sqrt{k_1! k_2! \dots k_m!}} \prod_{t=1}^m x_t^{k_t} \right]_{n=0, \dots, \infty; k_1+k_2+\dots+k_m=n}.$$

Since $k_1! k_2! \dots k_m! \rightarrow \infty$ for $n \rightarrow \infty$ the feature vector can be truncated at sufficiently large n .

Problem 3: Informally explain the difference between doing linear regression with a Gaussian RBF basis and doing linear regression with a Gaussian kernel.

When doing regression with a Gaussian RBF basis the centers of the RBFs can be placed anywhere while the Gaussian kernel implicitly uses the training samples as centers. Also, formally the centers of the Gaussian RBFs must be fixed in advance (before seeing the training samples), because the basis functions officially do not depend on the training samples. Linear regression using the Gaussian kernel can only be done in the dual representation, since the associated feature space has infinite dimensionality.

2 Kernel Perceptron

Problem 4: In this exercise, we develop a dual formulation of the perceptron algorithm. Using the perceptron learning rule you learned in the lecture, show that the learned weight vector can be written as a linear combination of the vectors $t^{(n)}\phi(\mathbf{x}^{(n)})$ where $t^{(n)} \in \{-1, +1\}$. Denote the coefficients of this linear combination by α_n and derive a formulation of the perceptron learning algorithm, and the predictive function for the perceptron in terms of the α_n . Show that the feature vector $\phi(\mathbf{x})$ enters only in the form of the kernel function $K(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x})^T \phi(\mathbf{y})$.

Perceptron:

$$\text{class } t^{(n)} = \text{sgn}(w^T \phi^{(n)}) \in \{-1, 1\}$$

correctly
classified $\phi^{(n)}$
($a^{(n)} = 0$)

$$\left\{ \begin{array}{l} \phi(x^{(n)}) = \phi^{(n)} \text{ in class } t^{(n)} = 1 \rightarrow w^T \phi^{(n)} > 0 \\ \phi(x^{(n)}) = \phi^{(n)} \text{ in class } t^{(n)} = -1 \rightarrow w^T \phi^{(n)} < 0 \\ \rightarrow w^T \phi^{(n)} t^{(n)} > 0 \end{array} \right.$$

incorrectly
classified $\phi^{(n)}$
($a^{(n)} = 1$)

$$\left\{ \begin{array}{l} \rightarrow w^T \phi^{(n)} t^{(n)} < 0 \end{array} \right.$$

→ minimize Error: $E = -\sum_n a^{(n)} w^T \phi^{(n)} t^{(n)}$

update rule: $w^{new} = w^{old} - \eta \nabla_w E$

$$w^{new} = \left\{ \begin{array}{ll} w^{old} - \eta a^{(n)} \phi^{(n)} t^{(n)} & \text{for online learning} \\ w^{old} - \eta \sum_n a^{(n)} \phi^{(n)} t^{(n)} & \text{for batch learning} \end{array} \right.$$

$$\xrightarrow{\eta = 1, w^{old} = 0} w = \sum_n a^{(n)} \phi^{(n)} t^{(n)}$$

$$\begin{aligned}
\text{class } t^{(t)} &= \text{sgn}(w^T \phi^{(t)}) \in \{-1, 1\} \\
&= \text{sgn} \left(\sum_n a^{(n)} t^{(n)} \phi^{(n)T} \phi^{(t)} \right) \\
&= \text{sgn} \left(\sum_n a^{(n)} t^{(n)} K(x^{(n)}, x^{(t)}) \right) = \text{sgn}(g(x^{(t)}))
\end{aligned}$$

→ learning rule for Kernel Perceptron:

Initialize $a_j = 0$ for all $j = 1, \dots, N$
do

for $t = 1, \dots, n$

$y' \leftarrow \text{sgn}(g(x^{(t)}))$

if $y' \neq t^{(t)}$ then $a_t \leftarrow a_t + 1$

end for

repeat until no classification mistakes

Thus the update rule $a_t \leftarrow a_t + 1$ for a misclassified training sample j corresponds to the update rule $w \leftarrow w + \phi(x^{(j)})$ of the conventional Perceptron training algorithm.

Part II: Problems w.r.t. Linear Classification

Problem 1: Given a set of data points $\{\mathbf{x}_n\}$, we can define the *convex hull* to be the set of all points \mathbf{x} given by

$$\mathbf{x} = \sum_n \alpha_n \mathbf{x}_n$$

where $\alpha_n \geq 0$ and $\sum_n \alpha_n = 1$. Consider a second set of points $\{\mathbf{y}_n\}$ together with their corresponding convex hull. By definition, the two sets of points will be linearly separable if there exists a vector \mathbf{w} and a scalar w_0 such that $\mathbf{w}^T \mathbf{x}_n + w_0 > 0$ for all \mathbf{x}_n , and $\mathbf{w}^T \mathbf{y}_n + w_0 < 0$ for all \mathbf{y}_n . Show that if their convex hulls intersect, the two sets of points cannot be linearly separable, and conversely that if they are linearly separable, their convex hulls do not intersect.

Problem 1: Given a set of data points $\{\mathbf{x}_n\}$, we can define the *convex hull* to be the set of all points \mathbf{x} given by

$$\mathbf{x} = \sum_n \alpha_n \mathbf{x}_n$$

where $\alpha_n \geq 0$ and $\sum_n \alpha_n = 1$. Consider a second set of points $\{\mathbf{y}_n\}$ together with their corresponding convex hull. By definition, the two sets of points will be linearly separable if there exists a vector \mathbf{w} and a scalar w_0 such that $\mathbf{w}^T \mathbf{x}_n + w_0 > 0$ for all \mathbf{x}_n , and $\mathbf{w}^T \mathbf{y}_n + w_0 < 0$ for all \mathbf{y}_n . Show that if their convex hulls intersect, the two sets of points cannot be linearly separable, and conversely that if they are linearly separable, their convex hulls do not intersect.

We prove *Convex hull intersect \Rightarrow not separable*. We have two sets of points, $\{\mathbf{x}_n\}$ and $\{\mathbf{y}_n\}$. Assume there is a \mathbf{z} that lies in both convex hulls, i.e. there are $\{\alpha_n\}$ and $\{\beta_n\}$ such that

$$\mathbf{z} = \sum_n \alpha_n \mathbf{x}_n = \sum_n \beta_n \mathbf{y}_n$$

Now, assume the two set of points are linearly separable, i.e. there is a \mathbf{w} for all $\{\mathbf{x}_n\}$ and $\{\mathbf{y}_n\}$ with $\mathbf{w}^T \mathbf{x}_n + w_0 > 0$ and $\mathbf{w}^T \mathbf{y}_n + w_0 < 0$. Now for $\mathbf{w}^T \mathbf{z}$ we get

$$\mathbf{w}^T \mathbf{z} = \sum_n \alpha_n \mathbf{w}^T \mathbf{x}_n > \sum_n \alpha_n (-w_0) = -w_0$$

but also

$$\mathbf{w}^T \mathbf{z} = \sum_n \beta_n \mathbf{w}^T \mathbf{y}_n < \sum_n \beta_n (-w_0) = -w_0$$

Logistic Regression (a **discriminative** classifier)

Problem 3 *Given a dataset $\mathcal{D} = \{(\mathbf{x}^n, c^n), n = 1, \dots, N\}$, where $c^n \in \{0, 1\}$, logistic regression uses the model $p(c = 1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b)$. Assuming the data is drawn independently and identically, show that the derivative of the log likelihood L with respect to \mathbf{w} is*

$$\nabla_{\mathbf{w}} L = \sum_{n=1}^N \left(c^n - \sigma(\mathbf{w}^T \mathbf{x}^n + b) \right) \mathbf{x}^n$$

Logistic Regression (a **discriminative** classifier)

Problem 3 *Given a dataset $\mathcal{D} = \{(\mathbf{x}^n, c^n), n = 1, \dots, N\}$, where $c^n \in \{0, 1\}$, logistic regression uses the model $p(c = 1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b)$. Assuming the data is drawn independently and identically, show that the derivative of the log likelihood L with respect to \mathbf{w} is*

$$\nabla_{\mathbf{w}} L = \sum_{n=1}^N \left(c^n - \sigma(\mathbf{w}^T \mathbf{x}^n + b) \right) \mathbf{x}^n$$

$$L = \sum_n c^n \log \sigma(\mathbf{w}^T \mathbf{x}^n + b) + (1 - c^n) \log(1 - \sigma(\mathbf{w}^T \mathbf{x}^n + b))$$

Using the derivative result $\sigma' = \sigma(1 - \sigma)$, and $\sigma_n \equiv \sigma(\mathbf{w}^T \mathbf{x}^n + b)$ we have

$$\nabla_{\mathbf{w}} L = \sum_{n=1} c^n (1 - \sigma_n) \mathbf{x}^n - (1 - c^n) \sigma_n \mathbf{x}^n = \sum_{n=1} (c^n - \sigma_n) \mathbf{x}^n$$

Multi-Class Logistic Regression

(a **discriminative** classifier)

Problem 4

Consider the softmax function for classifying an input vector \mathbf{x} into one of $c = 1, \dots, C$ classes using

$$p(c|\mathbf{x}) = \frac{e^{\mathbf{w}_c^T \mathbf{x}}}{\sum_{c'=1}^C e^{\mathbf{w}_{c'}^T \mathbf{x}}}$$

A set of input-class examples is given by $\mathcal{D} = \{(\mathbf{x}^n, c^n), n = 1, \dots, N\}$.

- 1. Relate this model to logistic regression when $C = 2$.*
- 2. Write down the log-likelihood L of the classes conditional on the inputs, assuming that the data is i.i.d.*

Multi-Class Logistic Regression

(a **discriminative** classifier)

Problem 4

Consider the softmax function for classifying an input vector \mathbf{x} into one of $c = 1, \dots, C$ classes using

$$p(c|\mathbf{x}) = \frac{e^{\mathbf{w}_c^T \mathbf{x}}}{\sum_{c'=1}^C e^{\mathbf{w}_{c'}^T \mathbf{x}}}$$

A set of input-class examples is given by $\mathcal{D} = \{(\mathbf{x}^n, c^n), n = 1, \dots, N\}$.

- 1. Relate this model to logistic regression when $C = 2$.*
- 2. Write down the log-likelihood L of the classes conditional on the inputs, assuming that the data is i.i.d.*

1. When $C = 2$ we have

$$p(c = 1|\mathbf{x}) = \frac{e^{\mathbf{w}_1^T \mathbf{x}}}{e^{\mathbf{w}_1^T \mathbf{x}} + e^{\mathbf{w}_2^T \mathbf{x}}} = \frac{1}{1 + e^{(\mathbf{w}_2 - \mathbf{w}_1)^T \mathbf{x}}} = \sigma(\mathbf{w}_2 - \mathbf{w}_1^T \mathbf{x})$$

For $C = 2$, this is therefore logistic regression parameterised by the difference $\mathbf{w} \equiv \mathbf{w}_2 - \mathbf{w}_1$. Note that there is therefore redundancy in the softmax parameterisation since only this difference plays a role in computing the probability. The same hold for $C > 2$ – by multiplying and dividing by $\exp(-\mathbf{w}_d^T \mathbf{x})$, for some chosen $d \in \{1, \dots, C\}$ only the differences $\mathbf{w}'_c \equiv \mathbf{w}_c - \mathbf{w}_d$ are required to specify the probability. This is potentially important to realise for optimisation purposes since there will therefore be equivalent solutions in the redundant parameterisation.

- 2.

$$L(\mathcal{D}) = \sum_n \left[\mathbf{w}_{c^n}^T \mathbf{x}^n - \log \sum_{d=1}^C e^{\mathbf{w}_d^T \mathbf{x}^n} \right]$$

Regularized Logistic Regression (MAP Solution to \mathbf{w})

(a **discriminative** classifier)

(Source: Jaaakkola.)

Problem 5

- a. Consider the data in Figure AAA, where we fit the model $p(y = 1|\mathbf{x}, \mathbf{w}) = \sigma(w_0 + w_1x_1 + w_2x_2)$. Suppose we fit the model by maximum likelihood, i.e., we minimize

$$J(\mathbf{w}) = -\ell(\mathbf{w}, \mathcal{D}_{\text{train}}) \quad (8.133)$$

where $\ell(\mathbf{w}, \mathcal{D}_{\text{train}})$ is the log likelihood on the training set. Sketch a possible decision boundary corresponding to $\hat{\mathbf{w}}$. (Copy the figure first (a rough sketch is enough), and then superimpose your answer on your copy, since you will need multiple versions of this figure). Is your answer (decision boundary) unique? How many classification errors does your method make on the training set?

- b. Now suppose we regularize only the w_0 parameter, i.e., we minimize

$$J_0(\mathbf{w}) = -\ell(\mathbf{w}, \mathcal{D}_{\text{train}}) + \lambda w_0^2 \quad (8.134)$$

Suppose λ is a very large number, so we regularize w_0 all the way to 0, but all other parameters are unregularized. Sketch a possible decision boundary. How many classification errors does your method make on the training set? Hint: consider the behavior of simple linear regression, $w_0 + w_1x_1 + w_2x_2$ when $x_1 = x_2 = 0$.

- c. Now suppose we heavily regularize only the w_1 parameter, i.e., we minimize

$$J_1(\mathbf{w}) = -\ell(\mathbf{w}, \mathcal{D}_{\text{train}}) + \lambda w_1^2 \quad (8.135)$$

Sketch a possible decision boundary. How many classification errors does your method make on the training set?

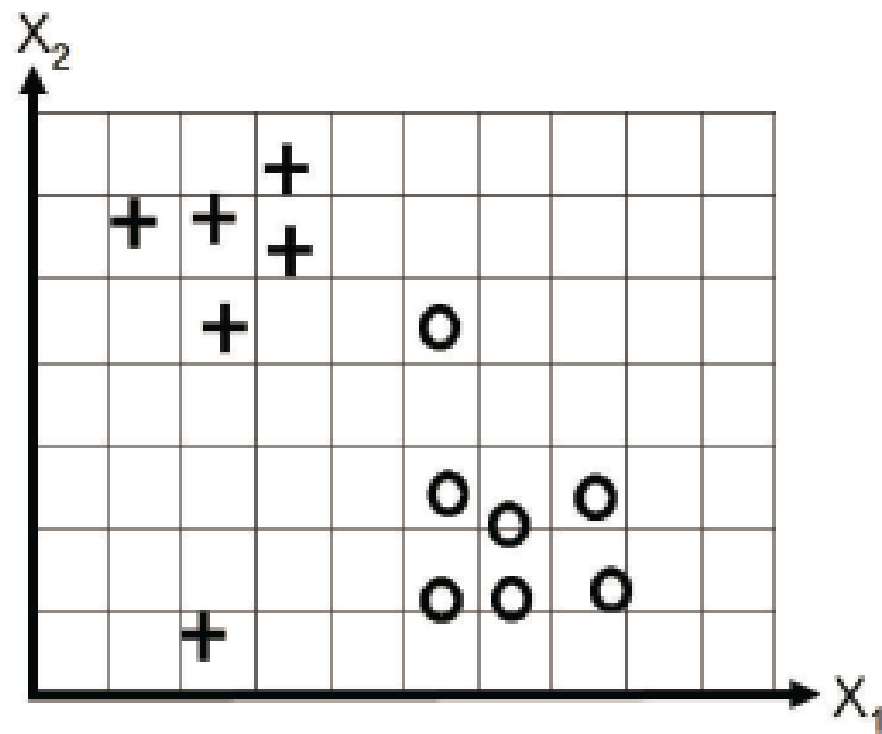


Figure AAA

1. Unregularized solution: see Figure 14.1(a). There is a unique ML decision boundary which makes 0 errors, but there are several possible decision boundaries which all makes **0 errors**.
2. Heavily regularizing w_0 sets $w_0 = 0$, so the line must go through the origin. There are several possible lines with different slopes, but all will make **1 error**. see Figure 14.1(b).
3. Regularizing w_1 makes the line horizontal since x_1 is ignored. This incurs **2 errors**. see Figure 14.1(c).
4. Regularizing w_2 makes the line vertical, since x_2 is ignored. This incurs **0 errors**. see Figure 14.1(d).

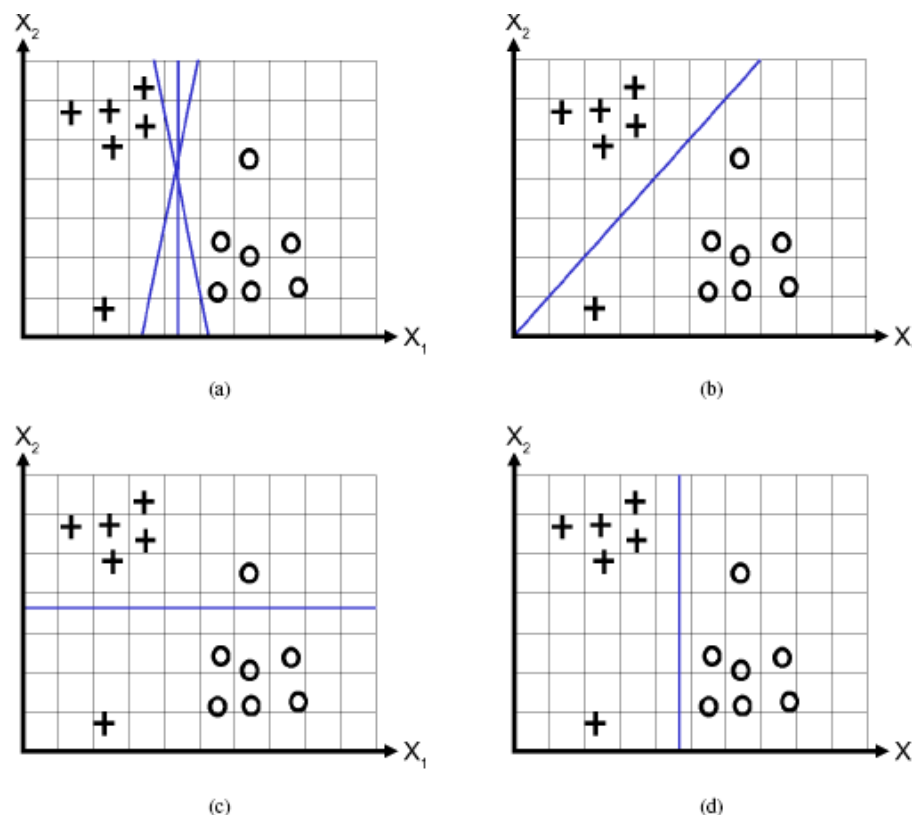


Figure 14.1: (a) Unregularized. (b) Regularizing w_0 . (c) Regularizing w_1 . (d) Regularizing w_2 .

Multi-Class Classification (Generative Multi-Class Classifiers)

Problem 6: Consider a generative classification model for K classes defined by prior class probabilities $p(y = k) = \pi_k$ and general class-conditional densities $p(\phi(x)|y = k, \theta_k)$ where $\phi(x)$ is the input feature vector and $\theta = \{\theta_k\}_{k=1}^K$ are further model parameters. Suppose we are given a training set $\mathcal{D} = \{(\phi(x^{(n)}), t^{(n)})\}_{n=1}^N = \{(\phi^{(n)}, t^{(n)})\}_{n=1}^N$ where $t^{(n)}$ is a binary target vector of length K that uses the 1-of- K (hot one) coding scheme, so that it has components $t_j^{(n)} = \delta_{jk}$ if pattern n is from class $y = k$. Assuming that the data points are iid, show that the maximum-likelihood solution for the prior probabilities is given by

$$\pi_k = \frac{N_k}{N}$$

where N_k is the number of data points assigned to class $y = k$.

The likelihood function of the parameters $\{\pi_k, \theta_k\}_{k=1}^K$ is given by

$$p(\mathcal{D}|\{\pi_k, \theta_k\}_{k=1}^K) = \prod_{n=1}^N \prod_{k=1}^K (p(\phi^{(n)}|\theta_k)\pi_k)^{t_k^{(n)}}$$

so the log-likelihood is given by

$$\log p(\mathcal{D}|\{\pi_k, \theta_k\}_{k=1}^K) = \sum_{n=1}^N \sum_{k=1}^K t_k^{(n)} \log \pi_k + \text{const in } \pi_k$$

In order to maximize the log likelihood with respect to π_k we need to preserve the constraint $\sum_k \pi_k = 1$. This can be done by introducing a Lagrange multiplier λ and maximizing

$$\sum_{n=1}^N \sum_{k=1}^K t_k^{(n)} \log \pi_k + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right).$$

Setting the derivative with respect to π_k equal to zero, we obtain

$$\lambda = \frac{1}{\pi_k} \sum_{n=1}^N t_k^{(n)} = \frac{1}{\pi_k} N_k.$$

Setting the derivative with respect to λ equal to zero, we obtain our constraint

$$\sum_{k=1}^K \pi_k = 1.$$

where we can now insert the previous result $\pi_k = \frac{N_k}{\lambda}$ and obtain

$$\lambda = \sum_k N_k = N.$$

So overall we obtain

$$\pi_k = \frac{N_k}{N}.$$

Linear Discriminant Analysis (LDA) (a *generative* classifier)

Problem 7

Topic:

Gaussian Discriminant Analysis (GDA)

(Variant (a): Quadratic Discriminant Analysis (QDA) : Σ_c class specific

Variant (b): Linear Discriminant Analysis (LDA): Σ shared among classes
)

vs. Logistic Regression

(compare lecture-slides on linear generative models)

Problem 7: Using the same classification model as in the previous question, now suppose that the class-conditional densities are given by Gaussian distributions with a shared covariance matrix, so that

$$p(\phi(x)|y = k, \theta_k) = p(\phi(x)|\theta_k) = \mathcal{N}(\phi \mid \mu_k, \Sigma).$$

Show that the maximum likelihood solution for the mean of the Gaussian distribution for class C_k is given by

$$\mu_k = \frac{1}{N_k} \sum_{\{n|\phi^{(n)} \in C_k\}} \phi^{(n)}$$

which represents the mean of those feature vectors assigned to class C_k .

Similarly, show that the maximum likelihood solution for the shared covariance matrix is given by

$$\Sigma = \sum_{k=1}^K \frac{N_k}{N} \mathbf{S}_k$$

where

$$\mathbf{S}_k = \frac{1}{N_k} \sum_{\{n|\phi^{(n)} \in C_k\}} (\phi^{(n)} - \mu_k)(\phi^{(n)} - \mu_k)^T.$$

Thus Σ is given by a weighted average of the covariances of the data associated with each class, in which the weighting coefficients N_k/N are the prior probabilities of the classes.

If we substitute $p(\phi|\theta_k) = \mathcal{N}(\phi | \mu_k, \Sigma)$ into $\log p(\mathcal{D}|\{\pi_k, \theta_k\}_{k=1}^K)$ and then use the definition of the multivariate Gaussian, we obtain

$$\log p(\mathcal{D}|\{\pi_k, \theta_k\}_{k=1}^K) = \frac{-1}{2} \sum_{n=1}^N \sum_{k=1}^K t_k^{(n)} (\log |\Sigma| + (\phi^{(n)} - \mu_k)^T \Sigma^{-1} (\phi^{(n)} - \mu_k) + \log \pi_k)$$

Dropping terms independent of μ_k and Σ we get

$$\log p(\mathcal{D}|\{\theta_k\}_{k=1}^K) = \frac{-1}{2} \sum_{n=1}^N \sum_{k=1}^K t_k^{(n)} (\log |\Sigma| + (\phi^{(n)} - \mu_k)^T \Sigma^{-1} (\phi^{(n)} - \mu_k))$$

Setting the derivative of the above equation w.r.t μ_k , (obtained using $\frac{\partial}{\partial x}(x^T a) = \frac{\partial}{\partial x}(a^T x) = a$), to zero we get

$$\sum_{n=1}^N t_k^{(n)} \Sigma^{-1} (\phi^{(n)} - \mu_k) = 0.$$

Making use of what we learned in the last problem, i.e.

$$\sum_{n=1}^N t_k^{(n)} = N_k,$$

we can re-arrange this to obtain

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N t_k^{(n)} \phi^{(n)}.$$

Using the trace trick ($a = \text{Tr}(a)$ for $a \in \mathbb{R}$ and $\text{Tr}(ABC) = \text{Tr}(BCA)$) we can rewrite our original expression $\log p(\mathcal{D}|\{\theta_k\}_{k=1}^K) = \frac{-1}{2} \sum_{n=1}^N \sum_{k=1}^K t_k^{(n)} (\log |\Sigma| + (\phi^{(n)} - \mu_k)^T \Sigma^{-1} (\phi^{(n)} - \mu_k))$ as

$$\log p(\mathcal{D}|\{\theta_k\}_{k=1}^K) = \frac{-1}{2} \sum_{n=1}^N \sum_{k=1}^K t_k^{(n)} (\log |\Sigma| + \text{Tr}(\Sigma^{-1} (\phi^{(n)} - \mu_k) (\phi^{(n)} - \mu_k)^T)).$$

We can now use $\frac{\partial}{\partial A} Tr(AB) = B^T$ and $\frac{\partial}{\partial A} \ln |A| = (A^{-1})^T$ and $|A^{-1}| = -|A|$ to calculate the derivative w.r.t. Σ^{-1} . Setting this to zero we obtain

$$\frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K t_k^{(n)} (\Sigma - (\phi^{(n)} - \mu_k) (\phi^{(n)} - \mu_k)^T) = 0.$$

Making use of $\sum_{n=1}^N t_k^{(n)} = N_k$, we can re-arrange this to obtain

$$\Sigma = \sum_{k=1}^K \frac{N_k}{N} \mathbf{S}_k$$

,

where

$$\mathbf{S}_k = \frac{1}{N_k} \sum_{n=1}^N t_k^{(n)} (\phi^{(n)} - \mu_k) (\phi^{(n)} - \mu_k)^T.$$

Note that we do not enforce that Σ should be symmetric, but the solution shows that it is automatically symmetric.

Questions to think about:

(a) What exactly is the difference between a generative classifier (e.g. GDA) and a discriminative classifier (e.g. Logistic Regression)?

(a) What are pros and cons of each approach?

(for „solution“ see Murphy 4.2 and 8.6)

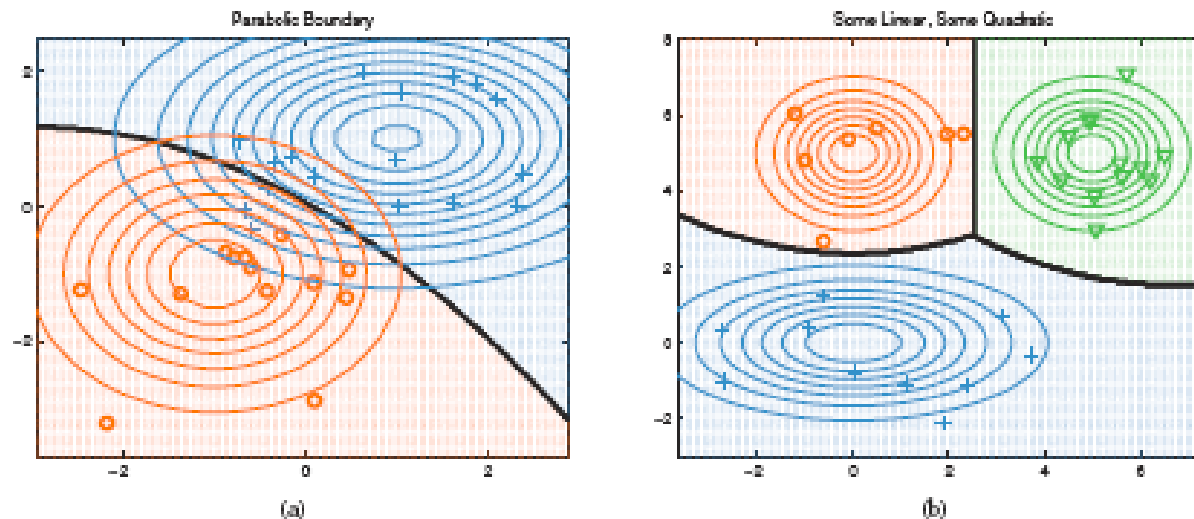


Figure 4.3 Quadratic decision boundaries in 2D for the 2 and 3 class case. Figure generated by `discrimAnalysisDboundariesDemo`.

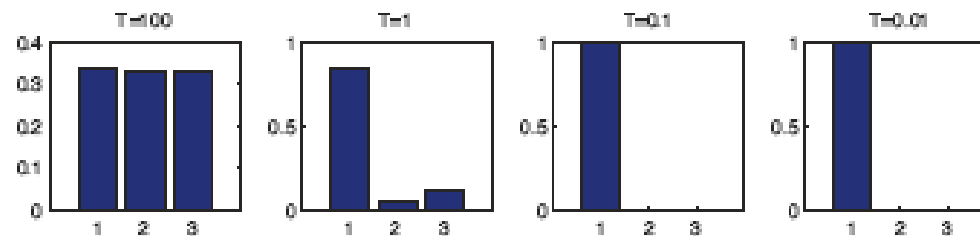


Figure 4.4 Softmax distribution $\mathcal{S}(\eta/T)$, where $\eta = (3, 0, 1)$, at different temperatures T . When the temperature is high (left), the distribution is uniform, whereas when the temperature is low (right), the distribution is "spiky", with all its mass on the largest element. Figure generated by `softmaxDemo2`.

$$p(y = c | \mathbf{x}, \theta) = \frac{\pi_c |2\pi \Sigma_c|^{-\frac{1}{2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \mu_c)^T \Sigma_c^{-1} (\mathbf{x} - \mu_c) \right]}{\sum_{c'} \pi_{c'} |2\pi \Sigma_{c'}|^{-\frac{1}{2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \mu_{c'})^T \Sigma_{c'}^{-1} (\mathbf{x} - \mu_{c'}) \right]}$$

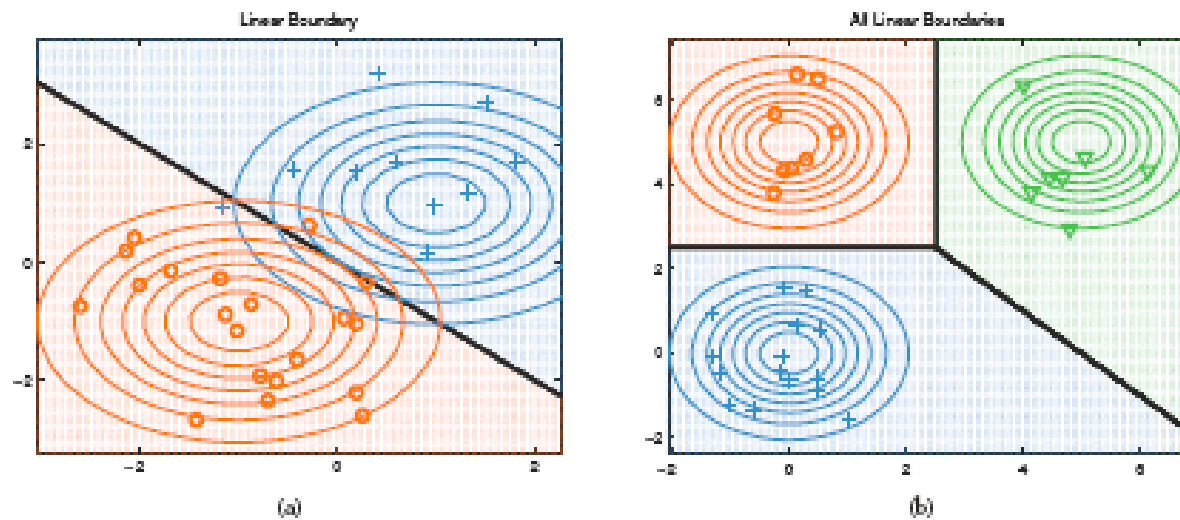


Figure 4.5 Linear decision boundaries in 2D for the 2 and 3 class case. Figure generated by `discrimAnalysisDboundariesDemo`.

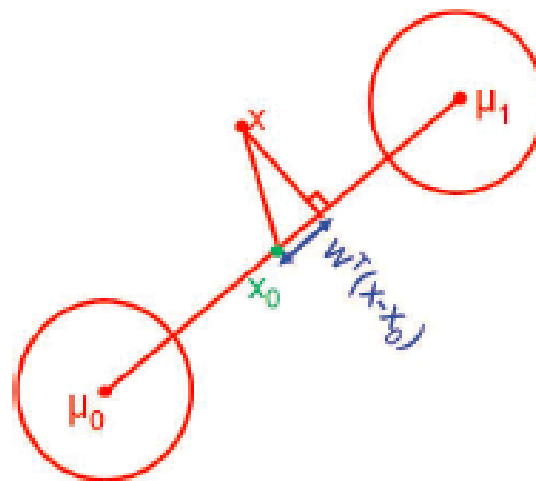
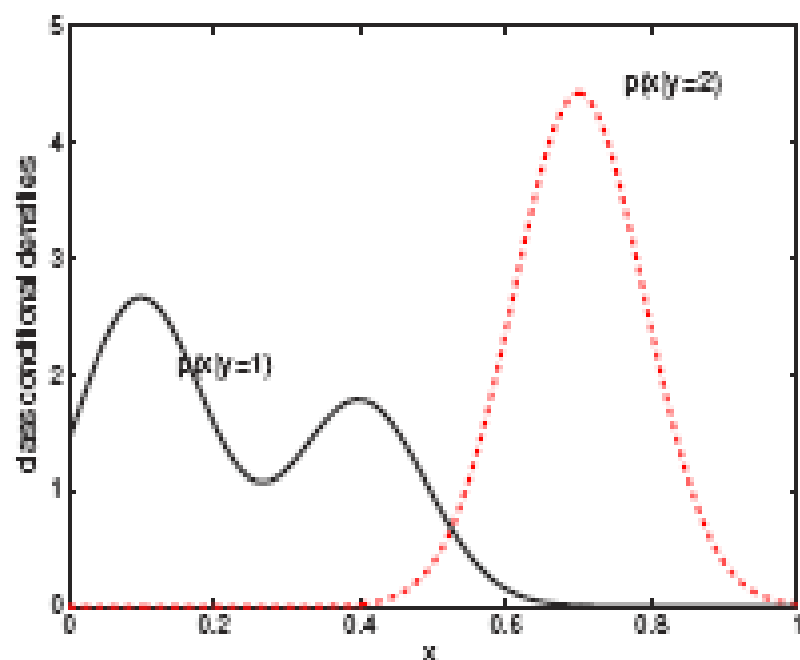
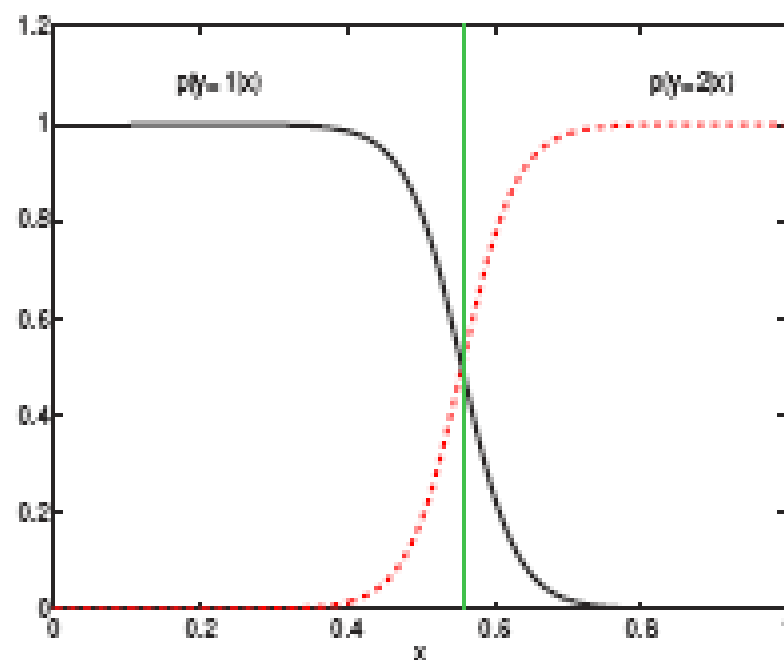


Figure 4.6 Geometry of LDA in the 2 class case where $\Sigma_1 = \Sigma_2 = I$.



(a)



(b)

Figure 8.10 The class-conditional densities $p(x|y = c)$ (left) may be more complex than the class posteriors $p(y = c|x)$ (right). Based on Figure 1.27 of (Bishop 2006a). Figure generated by generativeVsDiscrim.

Pros and cons of each approach

- **Easy to fit?** As we have seen, it is usually very easy to fit generative classifiers. For example, in Sections 3.5.1.1 and 4.2.4, we show that we can fit a naive Bayes model and an LDA model by simple counting and averaging. By contrast, logistic regression requires solving a convex optimization problem (see Section 8.3.4 for the details), which is much slower. g +1
- **Fit classes separately?** In a generative classifier, we estimate the parameters of each class conditional density independently, so we do not have to retrain the model when we add more classes. In contrast, in discriminative models, all the parameters interact, so the whole model must be retrained if we add a new class. (This is also the case if we train a generative model to maximize a discriminative objective Salojärvi et al. (2005).) g +1
- **Handle missing features easily?** Sometimes some of the inputs (components of \mathbf{x}) are not observed. In a generative classifier, there is a simple method for dealing with this, as we discuss in Section 8.6.2. However, in a discriminative classifier, there is no principled solution to this problem, since the model assumes that \mathbf{x} is always available to be conditioned on (although see (Madrina 2008) for some heuristic approaches). g +1
- **Can handle unlabeled training data?** There is much interest in semi-supervised learning, which uses unlabeled data to help solve a supervised task. This is fairly easy to do using generative models (see e.g., (Lasserre et al. 2006; Liang et al. 2007)), but is much harder to do with discriminative models. g +1
- **Symmetric in inputs and outputs?** We can run a generative model “backwards”, and infer probable inputs given the output by computing $p(\mathbf{x}|\mathbf{y})$. This is not possible with a discriminative model. The reason is that a generative model defines a joint distribution on \mathbf{x} and \mathbf{y} , and hence treats both inputs and outputs symmetrically. g +1
- **Can handle feature preprocessing?** A big advantage of discriminative methods is that they allow us to preprocess the input in arbitrary ways, e.g., we can replace \mathbf{x} with $\phi(\mathbf{x})$, which could be some basis function expansion, as illustrated in Figure 8.9. It is often hard to define a generative model on such pre-processed data, since the new features are correlated in complex ways. d +1
- **Well-calibrated probabilities?** Some generative models, such as naive Bayes, make strong independence assumptions which are often not valid. This can result in very extreme posterior class probabilities (very near 0 or 1). Discriminative models, such as logistic regression, are usually better calibrated in terms of their probability estimates. d +1

Exercise 4.20 Logistic regression vs LDA/QDA

(Source: Jaakkola.) Suppose we train the following binary classifiers via maximum likelihood.

- GaussI: A generative classifier, where the class conditional densities are Gaussian, with both covariance matrices set to \mathbf{I} (identity matrix), i.e., $p(\mathbf{x}|y = c) = \mathcal{N}(\mathbf{x}|\mu_c, \mathbf{I})$. We assume $p(y)$ is uniform.
- GaussX: as for GaussI, but the covariance matrices are unconstrained, i.e., $p(\mathbf{x}|y = c) = \mathcal{N}(\mathbf{x}|\mu_c, \Sigma_c)$.
- LinLog: A logistic regression model with linear features.
- QuadLog: A logistic regression model, using linear and quadratic features (i.e., polynomial basis function expansion of degree 2).

After training we compute the performance of each model M on the training set as follows:

$$L(M) = \frac{1}{n} \sum_{i=1}^n \log p(y_i | \mathbf{x}_i, \hat{\theta}, M) \quad (4.285)$$

(Note that this is the *conditional* log-likelihood $p(y|\mathbf{x}, \hat{\theta})$ and not the joint log-likelihood $p(y, \mathbf{x}|\hat{\theta})$.) We now want to compare the performance of each model. We will write $L(M) \leq L(M')$ if model M *must* have lower (or equal) log likelihood (on the training set) than M' , for any training set (in other words, M is worse than M' , at least as far as training set logprob is concerned). For each of the following model pairs, state whether $L(M) \leq L(M')$, $L(M) \geq L(M')$, or whether no such statement can be made (i.e., M might sometimes be better than M' and sometimes worse); also, for each question, briefly (1-2 sentences) explain why.

- GaussI, LinLog.
- GaussX, QuadLog.
- LinLog, QuadLog.
- GaussI, QuadLog.
- Now suppose we measure performance in terms of the average misclassification rate on the training set:

$$R(M) = \frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{y}(\mathbf{x}_i)) \quad (4.286)$$

Is it true in general that $L(M) > L(M')$ implies that $R(M) < R(M')$? Explain why or why not.

1. $\text{GaussI} \leq \text{LinLog}$. Both have logistic (sigmoid) posteriors $p(y|x, w) = \sigma(yw^T x)$, but LinLog is the logistic model which is trained to maximize $p(y|x, w)$. (GaussI may have high joint $p(y, x)$, but this does not necessarily mean $p(y|x)$ is high; LinLog can achieve the maximum of $p(y|x)$, so will necessarily do at least as well as GaussI.)
2. $\text{GaussX} \leq \text{QuadLog}$. Both have logistic posteriors with quadratic features, but QuadLog is the model of this class maximizing the average log probabilities.
3. $\text{LinLog} \leq \text{QuadLog}$. Logistic regression models with linear features are a subclass of logistic regression models with quadratic functions. The maximum from the superclass is at least as high as the maximum from the subclass.
4. $\text{GaussI} \leq \text{QuadLog}$. Follows from above inequalities.
5. Although one might expect that higher log likelihood results in better classification performance, in general, having higher average log $p(y|x)$ does not necessarily translate to higher or lower classification error. For example, consider linearly separable data. We have $L(\text{linLog}) > L(\text{GaussI})$, since maximum likelihood logistic regression will set the weights to infinity, to maximize the probability of the correct labels (hence $p(y_i|x_i, \hat{w}) = 1$ for all i). However, we have $R(\text{linLog}) = R(\text{gaussI})$, since the data is linearly separable. (The GaussI model may or may not set σ very small, resulting in possibly very large class conditional pdfs; however, the posterior over y is a discrete pmf, and can never exceed 1.)

As another example, suppose the true label is always 1 (as opposed to 0), but model M always predicts $p(y = 1|x, M) = 0.49$. It will always misclassify, but it is at least close to the decision boundary. By contrast, there might be another model M' that predicts $p(y = 1|x, M') = 1$ on even-numbered inputs, and $p(y = 1|x, M') = 0$ on odd-numbered inputs. Clearly $R(M') = 0.5 < R(M) = 1$, but $L(M') = -\infty < L(M) = \log(0.49)$.

Hinge loss

The hinge loss is given as

$$\mathcal{L}(\mathbf{x}_i) = \max(0, 1 - y_i \tilde{z}_i),$$

where $y_i = \mathbf{x}_i^T \mathbf{w}$ is the model output and z_i the target variable ($w_0 = b$ and $x_{i,0} = 1$).

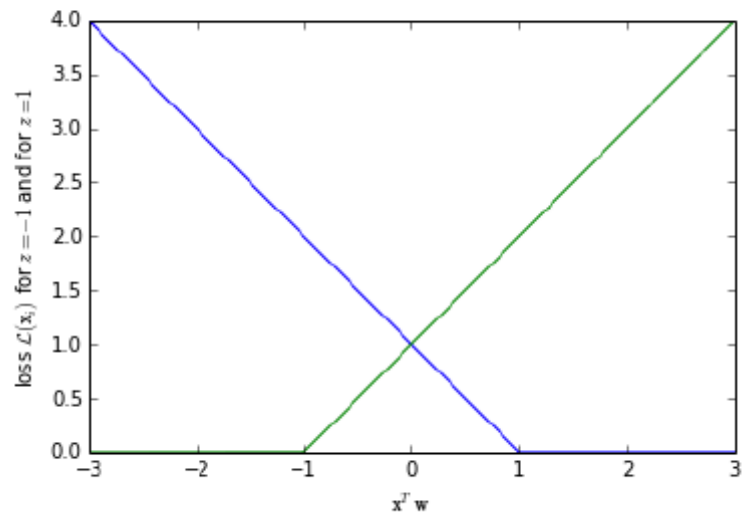
Note that in this case the computation uses class labels $\tilde{z} = 2z - 1 \in \{-1, 1\}$ instead of $z \in \{0, 1\}$.

For multiple samples \mathbf{X} and respective outputs \mathbf{y} and $\tilde{\mathbf{z}}$ the loss is $\mathcal{L}(\mathbf{X}) = \sum_i \mathcal{L}(\mathbf{x}_i)$.

Problem 8 Try to understand what the hinge loss does and explain it in a few of words

Problem 9 Derive a sub-gradient of the hinge loss $\frac{d\mathcal{L}(\mathbf{x}_i)}{d\mathbf{w}}$. (You may assume that $w_0 = b$. Take 0 for the sub-gradient at non-differentiable points. Otherwise deriving a sub-gradient is equivalent to deriving the gradient.)

The decision boundary of the hinge loss is where $y = \mathbf{x}^T \mathbf{w} = 0$. The loss function for $z = 1$ and $z = -1$ looks like:



The hinge loss attempts to create a small margin between the decision boundary and the closest training samples. It slightly penalizes even samples that are on the correct side of the decision boundary but are within the margin ($|\mathbf{x}^T \mathbf{w}| < 1$) and penalizes incorrectly classified samples linearly with their distance from the (inner) margin.

The hinge loss is not differentiable at $y_i \tilde{z}_i = 1$. However, since the hinge loss is convex, we can derive a subgradient:

$$\begin{aligned} \frac{d\mathcal{L}(x_i)}{dw} &= \begin{cases} \frac{d(1-y_i \tilde{z}_i)}{dw} & \text{if } y_i \tilde{z}_i < 1 \\ 0 & \text{else} \end{cases} \\ &= \begin{cases} \frac{d(1-\tilde{z}_i x_i^T w)}{dw} & \text{if } y_i \tilde{z}_i < 1 \\ 0 & \text{else} \end{cases} \\ &= \begin{cases} -\tilde{z}_i x_i & \text{if } y_i \tilde{z}_i < 1 \\ 0 & \text{else} \end{cases} \end{aligned}$$

For multiple outputs $y = Xw$ and targets \tilde{z} :

$$\frac{d\mathcal{L}(X)}{dw} = \sum_i -z_i x_i, \quad i \in \{1, \dots, N\} : y_i z_i < 1$$

Soft Zero-one loss

The soft zero-one loss is given as

$$\mathcal{L}(\mathbf{x}_i) = (\sigma(\beta y_i) - z_i)^2,$$

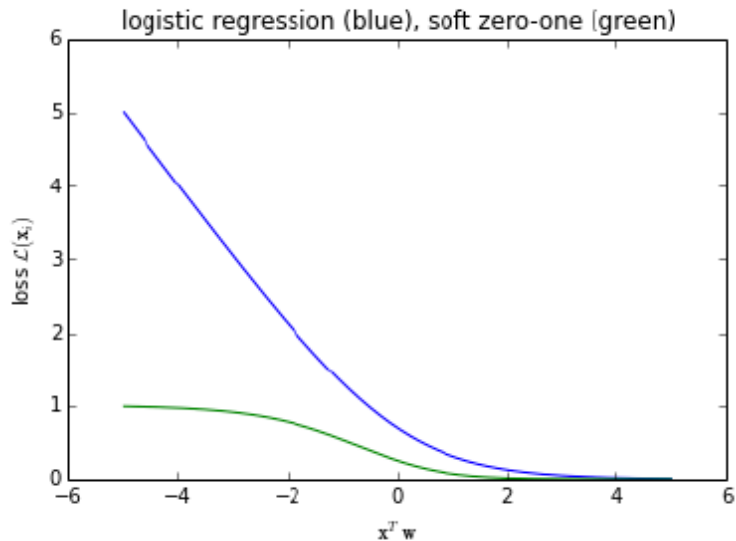
with $y_i = \mathbf{x}_i^T \mathbf{w}$ the model output and z_i the target variable.

Problem 10 Explain the soft zero-one loss in a few words.

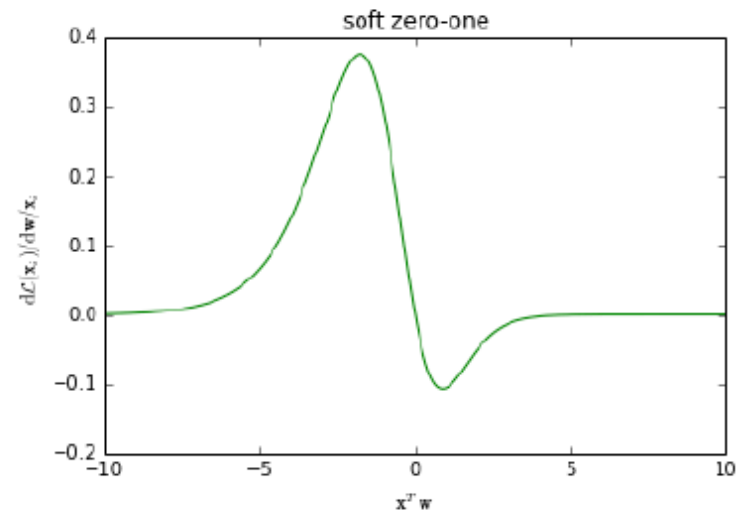
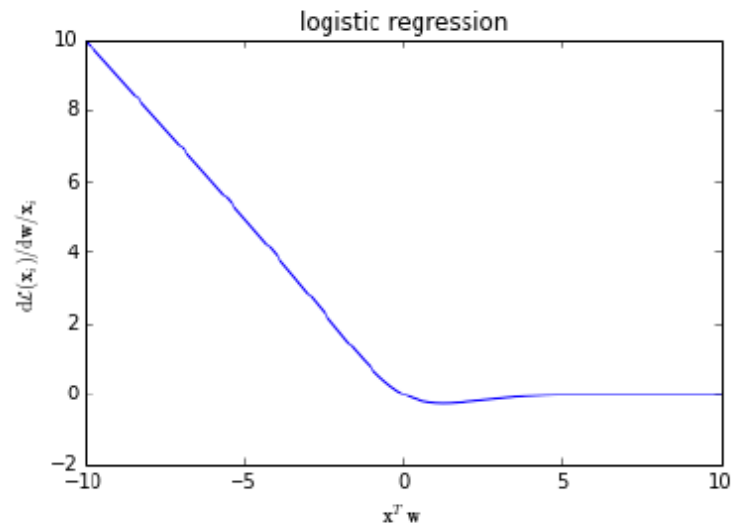
Problem 11 Derive the gradient $\frac{d\mathcal{L}(\mathbf{x}_i)}{d\mathbf{w}}$.

It is difficult to see the difference between the loss of logistic regression and the soft zero-one loss. One difference is the use of the β parameter which controls the steepness of the sigmoid function. However, a steeper or flatter sigmoid can also be attained by increasing or decreasing the magnitude of the weight vector.

The main difference, however, is in the loss function. The logistic regression loss (the negative log likelihood) is unbounded. A sample that is incorrectly labeled and far from the decision boundary has a loss that increases linearly with the distance from the boundary. In contrast, the soft zero-one loss is bounded. This is the reason why it is more robust against outliers because even incorrectly labeled samples that are far from the decision boundary ($|x^T w| \gg 0$) will change the loss by no more than 1.



Looking at the gradient can give us additional insight. Interestingly, the gradient of the soft zero-one loss approaches zero for incorrectly classified samples that are far from the decision boundary. Therefore, they almost do not contribute to the gradient, whereas in logistic regression their impact is linear with the distance from the decision boundary:



$$\begin{aligned} \frac{d\mathcal{L}(\mathbf{x}_i)}{d\mathbf{w}} &= \frac{d(\sigma(\beta y_i) - z_i)^2}{d\mathbf{w}} \\ &= 2(\sigma(\beta y_i) - z_i) \sigma(\beta y_i) (1 - \sigma(\beta y_i)) \beta \mathbf{x}_i \end{aligned}$$