

1: AML Pg. 79, Ex. 3.1**Answer:**

PLA will never stop updating if the data is not linearly separable because PLA only stops on the condition that there are no misclassified datapoints. If, however, the data is not linearly separable, any PLA-generated decision boundary, by being linear, will result in at least one misclassified datapoint.

2: AML Pg. 99, Ex. 3.7**Answer:** a) Showing

$$-\frac{1}{N} \sum_{n=1}^N \frac{y_n \mathbf{x}_n}{1 + \exp(y_n \mathbf{w}^T \mathbf{x}_n)} = \frac{1}{N} \sum_{n=1}^N -y_n \mathbf{x}_n \theta(-y_n \mathbf{w}^T \mathbf{x}_n)$$

is the same as showing

$$\frac{1}{1 + \exp(y_n \mathbf{w}^T \mathbf{x}_n)} = \theta(-y_n \mathbf{w}^T \mathbf{x}_n) = \frac{\exp(-y_n \mathbf{w}^T \mathbf{x}_n)}{1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n)}$$

Define $\alpha = y_n \mathbf{w}^T \mathbf{x}_n$

$$\frac{1}{1 + \exp \alpha} * \frac{\exp -\alpha}{\exp -\alpha} = \frac{\exp -\alpha}{1 + \exp -\alpha}$$

Substituting in α ...

$$\frac{1}{1 + \exp(y_n \mathbf{w}^T \mathbf{x}_n)} = \frac{\exp(-y_n \mathbf{w}^T \mathbf{x}_n)}{1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n)}$$

b) If an example is misclassified, α is negative, so $\exp \alpha < 1$. The denominator of the error gradient of this example is therefore smaller than if the example were correctly classified, meaning the gradient is larger. A misclassified example therefore contributed more to a change in weights.

3: AML Pg. 94, Ex. 3.8**Answer:**

The unit vector $\hat{\mathbf{v}}$ only gives the largest E_{in} decrease for small η because $\hat{\mathbf{v}}$ is defined by $\hat{\mathbf{v}} \propto \nabla E_{in}$. Large η will adjust weights to be far from these initial weights, at which point the unit vector measuring ∇E_{in} slope will have changed from $\hat{\mathbf{v}}$. These new weights may therefore overshoot the minimum E_{in} . a)

$$h(\mathbf{x}(t)) = \mathbf{w}^T(t) \mathbf{x}(t)$$

For misclassification, $y(t) = -h(\mathbf{x}(t))$.*Proof by cases:*

Case 1:

If $y(t) = +1$, $\mathbf{w}^T(t) \mathbf{x}(t) = -1$, so $y(t) \mathbf{w}^T(t) \mathbf{x}(t) = -1 < 0$

Case 2:

If $y(t) = -1$, $\mathbf{w}^T(t)\mathbf{x}(t) = +1$, so $y(t)\mathbf{w}^T(t)\mathbf{x}(t) = -1 < 0$

b) By eq. 1.3:

$$\begin{aligned} y(t)\mathbf{w}^T(t+1)\mathbf{x}(t) &= y(t)[\mathbf{w}(t) + y(t)\mathbf{x}(t)]^T\mathbf{x}(t) \\ &= y(t)\mathbf{w}^T(t)\mathbf{x}(t) + y^2(t)\mathbf{x}^T(t)\mathbf{x}(t) \\ \therefore y(t)\mathbf{w}^T(t+1)\mathbf{x}(t) &> y(t)\mathbf{w}^T(t)\mathbf{x}(t) \end{aligned}$$

holds if and only if

$$y^2(t)\mathbf{x}^T(t)\mathbf{x}(t) > 0$$

Both $y^2(t)$ and $\mathbf{x}^T(t)\mathbf{x}(t)$ are positive – the former because any scalar squared is positive, and the latter because $\mathbf{x}^T(t)\mathbf{x}(t)$ is a sum of the squared scalars within $\mathbf{x}(t)$. Both $y^2(t)\mathbf{x}^T(t)\mathbf{x}(t) > 0$ and subsequently $y(t)\mathbf{w}^T(t+1)\mathbf{x}(t) > y(t)\mathbf{w}^T(t)\mathbf{x}(t)$ are therefore true.

c) If $h(\mathbf{x}(t)) = -1$ and the point $\mathbf{x}(t)$ was misclassified, $y(t) = +1$. The update rule $\mathbf{w}(t+1) = \mathbf{w}(t) + y(t)\mathbf{x}(t)$ will therefore make the new decision boundary more positive. That is, for any t , $\mathbf{w}^T(t+1)\mathbf{x}(t)$ is more likely to be positive than $\mathbf{w}^T(t)\mathbf{x}(t)$. A future point similar to $\mathbf{x}(t)$ is therefore less likely to be misclassified as negative. The same argument holds with $h(\mathbf{x}(t)) = +1$ and $y(t) = -1$, where in this case the decision boundary is made more negative. Each misclassified point causes an adjustment that makes nearby points less likely to be misclassified in the same way.

4: AML Pg. 97, Ex. 3.9

Answer:

a)

See code notebook

b)

As one can see in the notebook graph, $e_{class}(s, y)$ makes its closest approach to $e_{sq}(s, y)$ when $s = 0$ and when $s = 1$. In the former case, $\lim_{s \rightarrow 0^+} e_{class}(s, 1) = 1$, and $e_{sq}(0, 1) = 1$. For $s < 0$, $e_{class} = 1$ and e_{sq} is decreasing, so for all $s \leq 0$, $e_{class} \leq e_{sq}$.

At $s = 1$, $e_{class}(1, 1) = 0$ and $e_{sq}(1, 1) = 1$. For $0 < s < 1$, e_{sq} is decreasing and $e_{class} = 0$, so $e_{class} < e_{sq}$. Lastly, for $s > 1$, $e_{class} = 0$ and e_{sq} is increasing from 0.

For all s , therefore, $e_{class} < e_{sq}$.

c)

Unlike $e_{sq}(s, 1)$, $e_{log}(s, 1)$ monotonically decreases with s . This is because

$\frac{de_{log}(s, 1)}{ds} = \frac{1}{1 + \exp(-ys)} * \exp(-ys) * (-y)$, which is negative for all s . This monotonicity is also evident in the graph. Indeed, one can infer from the graph that $\lim_{s \rightarrow \infty} e_{log}(s, 1) = 0$. This holds mathematically as,

$$\lim_{s \rightarrow \infty} e_{log}(s, 1) = \lim_{s \rightarrow \infty} \ln(1 + \exp(-ys)) = \ln(1 + \exp(-\infty)) = \ln(1) = 0$$

5: PRML Pg. 134, Ex. 2.43

Answer: First, note that $\exp(\frac{-|x|^q}{2\sigma^2})$ is even, such that

$$\int_0^\infty \exp(\frac{-x^q}{2\sigma^2})dx = \frac{1}{2} \int_{-\infty}^\infty \exp(\frac{-|x|^q}{2\sigma^2})dx$$

So...

$$\frac{q}{2(2\sigma^2)^{1/q}\Gamma(\frac{1}{q})} \int_{-\infty}^\infty \exp(\frac{-|x|^q}{2\sigma^2})dx = \frac{q}{(2\sigma^2)^{1/q}\Gamma(\frac{1}{q})} \int_0^\infty \exp(\frac{-x^q}{2\sigma^2})dx$$

If we define $u = \frac{x^q}{2\sigma^2}$, then substituting in we get

$$\begin{aligned} \frac{q}{(2\sigma^2)^{1/q}\Gamma(\frac{1}{q})} \int_0^\infty \exp(\frac{-x^q}{2\sigma^2})dx &= \frac{q}{(2\sigma^2)^{1/q}\Gamma(\frac{1}{q})} \int_0^\infty \frac{(2\sigma^2)^{1/q}}{q} u^{1/q} \exp(-u)du \\ &= (\frac{q}{(2\sigma^2)^{1/q}\Gamma(\frac{1}{q})}) * (\frac{(2\sigma^2)^{1/q}}{q} \Gamma(\frac{1}{q})) = 1 \end{aligned}$$

As the integral evaluates to one, this probability distribution is normalized.

For $q=2$, we get

$$p(x|\sigma^2, 2) = \frac{1}{(\sigma^2)^{1/2}\Gamma(\frac{1}{2})} \exp(-\frac{|x|^2}{2\sigma^2})$$

and since $\Gamma(1/2) = \sqrt{2\pi}...$

$$p(x|\sigma^2, 2) = \frac{1}{\sigma\sqrt{2\pi}} \exp(-\frac{|x|^2}{2\sigma^2})$$

This is the Gaussian distribution.

The target variable can be expressed as $t = y(\mathbf{x}, \mathbf{w}) + \epsilon$, and therefore $-\epsilon$, a quantity we would like to maximize, is $y(\mathbf{x}, \mathbf{w}) - t$. The probability of achieving a target outcome is the sum of log likelihoods, where likelihood is the probability of $-\epsilon$. We therefore have

$$\begin{aligned} \ln p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) &= \sum_{n=1}^N \ln p(y(\mathbf{x}_n, \mathbf{w}) - t_n | \sigma^2, q) \\ &= \frac{-1}{2\sigma^2} \sum_{n=1}^N |y(\mathbf{x}_n, \mathbf{w}) - t_n|^q + N \ln \frac{q}{2(2\sigma^2)^{\frac{1}{q}}\Gamma(\frac{1}{q})} \\ &= \frac{-1}{2\sigma^2} \sum_{n=1}^N |y(\mathbf{x}_n, \mathbf{w}) - t_n|^q - \frac{N \ln 2\sigma^2}{q} + \end{aligned}$$

6: PRML Pg. 173, Ex. 3.1

Answer:

We know that $\tanh a = \frac{e^a - e^{-a}}{e^a + e^{-a}}$ Let's start from the following:

$$2\sigma(2a) - 1 = \frac{2}{1 + e^{-2a}} - 1$$

$$\begin{aligned}
&= \frac{2}{1+e^{-2a}} * \frac{e^{2a}}{e^{2a}} - 1 = \frac{2e^{2a}}{1+e^{2a}} - 1 \\
&= \frac{2e^{2a}}{1+e^{2a}} - \frac{1+e^{2a}}{1+e^{2a}} = \frac{2e^{2a} - (1+e^{2a})}{1+e^{2a}} \\
&= \frac{e^{2a} - 1}{1+e^{2a}} = \frac{e^{2a} - 1}{1+e^{2a}} * \frac{e^{-a}}{e^{-a}} \\
&= \frac{e^a - e^{-a}}{e^a + e^{-a}}
\end{aligned}$$

That is, $2\sigma(2a) - 1 = \tanh a$.

Now start with $y(x, \mathbf{w}) = w_0 + \sum_{j=1}^M \mathbf{w}_j \sigma(2a_j)$ where $a_j = \frac{x - \mu_j}{2s}$

$$y(x, \mathbf{w}) = w_0 + \sum_{j=1}^M \frac{1}{2} \mathbf{w}_j (2\sigma(2a_j) + 1 - 1)$$

By the result $2\sigma(2a) - 1 = \tanh a$...

$$y(x, \mathbf{w}) = w_0 + \sum_{j=1}^M \frac{1}{2} \mathbf{w}_j (\tanh a_j + 1)$$

Defining $u_j = \frac{\mathbf{w}_j}{2}$...

$$\begin{aligned}
&= w_0 + \sum_{j=1}^M u_j \tanh a_j + \sum_{j=1}^M \frac{\mathbf{w}_j}{2} \\
&= u_0 + \sum_{j=1}^M u_j \tanh a_j
\end{aligned}$$

where I define $u_0 = w_0 + \sum_{j=1}^M \frac{\mathbf{w}_j}{2}$.

7: PRML Pg. 174, Ex. 3.2

Answer: Suppose we have a vector \mathbf{v} that can be expressed by a sum of its orthogonal components as $\mathbf{v} = \mathbf{v}_1 + \mathbf{v}_2$, where $\mathbf{v}_1 \in \Phi$ and $\mathbf{v}_2 \in \Phi^\perp$. There then exists a vector \mathbf{w} such that $\Phi \mathbf{w} = \mathbf{v}_1$. Let's operate on \mathbf{v}_1 as follows:

$$\Phi(\Phi^T \Phi)^{-1} \Phi^T \mathbf{v}_1 = \Phi(\Phi^T \Phi)^{-1} \Phi^T \Phi \mathbf{w}$$

But we know that $(\Phi^T \Phi)^{-1} \Phi^T \Phi = 1$, so

$$\begin{aligned}
\Phi(\Phi^T \Phi)^{-1} \Phi^T \mathbf{v}_1 &= \Phi \mathbf{w} \\
&= \mathbf{v}_1
\end{aligned}$$

The operator $\Phi(\Phi^T \Phi)^{-1} \Phi^T$, then, returns \mathbf{v}_1 when it operates on \mathbf{v}_1 .

By contrast, take \mathbf{v}_2 , for which we know (by the definition of orthogonality) that $\Phi^T \mathbf{v}_2 = 0$. Therefore...

$$\Phi(\Phi^T \Phi)^{-1} \Phi^T \mathbf{v}_2 = \Phi(\Phi^T \Phi)^{-1} 0 = 0$$

In conclusion, we can take any vector \mathbf{v} and operate on it as follows:

$$\begin{aligned}\Phi(\Phi^T\Phi)^{-1}\Phi^T\mathbf{v} &= \Phi(\Phi^T\Phi)^{-1}\Phi^T(\mathbf{v}_1 + \mathbf{v}_2) \\ &= \Phi(\Phi^T\Phi)^{-1}\Phi^T\mathbf{v}_1 \\ &= \mathbf{v}_1\end{aligned}$$

The vector \mathbf{v} , therefore, is projected onto the columns of Φ by $\Phi(\Phi^T\Phi)^{-1}\Phi^T$, as $\mathbf{v}_1 \in \Phi$.

If we consider Φ to be our matrix of features and the vector \mathbf{v} to be our target values \mathbf{t} , then $\Phi(\Phi^T\Phi)^{-1}\Phi^T\mathbf{t}$ projects \mathbf{t} onto the columns of our features, which is a manifold \mathcal{S} of dimensionality M (the number of features).

8: PRML Pg. 220, Ex. 4.1

Answer:

We can define an intersection of convex holds to be the point $\mathbf{z} = \sum_n \alpha_n \mathbf{x}_n = \sum_n \alpha_n \mathbf{y}_n$ for two sets of points $\{\mathbf{x}_n\}$ and $\{\mathbf{y}_n\}$.

For separability, we require a \mathbf{w}^T and w_0 such that $\mathbf{w}^T \mathbf{x}_n + w_0 > 0$ for all n and $\mathbf{w}^T \mathbf{y}_n + w_0 < 0$ for all n . Consider, however, the point of intersection...

$$\mathbf{w}^T \mathbf{z} + w_0 = \mathbf{w}^T \sum_n \alpha_n \mathbf{x}_n + w_0$$

Since, by the definition of α_n , $\sum_n \alpha_n = 1$...

$$\mathbf{w}^T \mathbf{z} + w_0 = \mathbf{w}^T \sum_n \alpha_n \mathbf{x}_n + \sum_n \alpha_n w_0$$

Rearranging and acknowledging that \mathbf{w}^T is independent of n ...

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

However, by the same procedure, we also find the following:

$$\mathbf{w}^T \mathbf{z} + w_0 = \sum_n \alpha_n (\mathbf{w}^T \mathbf{y}_n + w_0)$$

We have found that $\sum_n \alpha_n (\mathbf{w}^T \mathbf{x}_n + w_0) = \sum_n \alpha_n (\mathbf{w}^T \mathbf{y}_n + w_0)$. Therefore, if $\mathbf{w}^T \mathbf{x}_n + w_0$ is positive, negative, or zero, $\mathbf{w}^T \mathbf{y}_n + w_0$ also is. This means $\{\mathbf{x}_n\}$ and $\{\mathbf{y}_n\}$ are not linearly separable, by definition.

Assuming linear separability gives us the following:

$$\mathbf{w}^T \mathbf{x}_n + w_0 > 0$$

and

$$\mathbf{w}^T \mathbf{y}_n + w_0 < 0$$

or vice versa. From before, we know that

$$\begin{aligned}\mathbf{w}^T \mathbf{z}_1 + w_0 &= \sum_n \alpha_n (\mathbf{w}^T \mathbf{x}_n + w_0) \\ \mathbf{w}^T \mathbf{z}_2 + w_0 &= \sum_n \alpha_n (\mathbf{w}^T \mathbf{y}_n + w_0)\end{aligned}$$

where

$$\begin{aligned}\mathbf{z}_1 &= \sum_n \alpha_n \mathbf{x}_n \\ \mathbf{z}_2 &= \sum_n \alpha_n \mathbf{y}_n\end{aligned}$$

This means that

$$\begin{aligned}\mathbf{w}^T \mathbf{z}_1 + w_0 &\neq \mathbf{w}^T \mathbf{z}_2 + w_0 \\ \mathbf{z}_1 &\neq \mathbf{z}_2\end{aligned}$$

Since \mathbf{z}_1 and \mathbf{z}_2 are the convex hulls of \mathbf{x}_n and \mathbf{y}_n , respectively, the hulls do not intersect.

9: PRML Pg. 221, Ex. 4.7

Answer:

Show $\sigma(-a) = 1 - \sigma(a)$.

$$\begin{aligned}\sigma(-a) &= \frac{1}{1 + e^a} \\ &= \frac{e^{-a}}{1 + e^{-a}} \\ &= \frac{1 + e^{-a} - 1}{1 + e^{-a}} \\ &= 1 - \frac{1}{1 + e^{-a}} \\ &= 1 - \sigma(a)\end{aligned}$$

Show $\sigma^{-1}(y) = \ln \left[\frac{y}{1-y} \right]$

$$\begin{aligned}\sigma^{-1}(y) \rightarrow y &= \sigma(a) = \frac{1}{1 + e^{-a}} \\ 1 + e^{-a} &= \frac{1}{y} \\ \sigma^{-1}(y) = a &= -\ln \left[\frac{1}{y} - 1 \right]\end{aligned}$$

By the rules of logarithms...

$$\begin{aligned}\sigma^{-1}(y) &= \ln \left[\frac{1}{y} - 1 \right]^{-1} \\ &= \ln \left[\frac{1}{\frac{1}{y} - 1} \right]\end{aligned}$$

$$\begin{aligned}
&= \ln \left[\frac{1}{\frac{1}{y} - 1} * \frac{y}{y} \right] \\
&= \ln \left[\frac{y}{1 - y} \right]
\end{aligned}$$

10: PRML Pg. 222, Ex. 4.12

Answer:

$$\begin{aligned}
\frac{d\sigma}{da} &= \frac{d}{da} \left(\frac{1}{1 + e^{-a}} \right) \\
&= \frac{d}{da} (1 + e^{-a})^{-1} \\
&= -1 * (1 + e^{-a})^{-2} * e^{-a} * -1 \\
&= \frac{e^{-a}}{(1 + e^{-a})^2} \\
&= \frac{1}{1 + e^{-a}} * \frac{e^{-a}}{1 + e^{-a}} * \frac{e^a}{e^a} \\
&= \sigma(a) * \frac{1}{1 + e^a} \\
&= \sigma(a) * \sigma(-a)
\end{aligned}$$

And by the result of PRML Ex. 4.7...

$$\frac{d\sigma}{da} = \sigma(a)(1 - \sigma(a))$$

11: PRML Pg. 222, Ex. 4.13

Answer:

$$E(\mathbf{w}) = - \sum_{n=1}^N [t_n \ln y_n + (1 - t_n) \ln(1 - y_n)]$$

where $y_n = \sigma(\mathbf{w}^T \phi_n)$ We can distribute the gradient into the sum as follows:

$$\nabla_w E = - \sum_{n=1}^N [t_n \nabla_w \ln y_n + (1 - t_n) \nabla_w \ln(1 - y_n)]$$

since the target is independent of the weights. By the result of PRML Ex. 10 and the chain rule...

$$\nabla_w E = - \sum_{n=1}^N \left[t_n \left(\frac{1}{\sigma(\mathbf{w}^T \phi_n)} * \sigma(\mathbf{w}^T \phi_n) * (1 - \sigma(\mathbf{w}^T \phi_n)) * \phi_n \right) + (1 - t_n) \left(\frac{1}{1 - \sigma(\mathbf{w}^T \phi_n)} * -\sigma(\mathbf{w}^T \phi_n) * (1 - \sigma(\mathbf{w}^T \phi_n)) * \phi_n \right) \right]$$

Distributing...

$$\nabla_w E = - \sum_{n=1}^N [(t_n \phi_n - t_n \phi_n \sigma(\mathbf{w}^T \phi_n)) + (-\phi_n \sigma(\mathbf{w}^T \phi_n) + t_n \phi_n \sigma(\mathbf{w}^T \phi_n))]$$

The second and fourth terms cancel, leaving the following:

$$\nabla_w E = - \sum_{n=1}^N \phi_n (t_n - \sigma(\mathbf{w}^T \phi_n))$$

12: PRML Pg. 222, Ex. 4.14

Answer:

Assume there are two classes C_1 and C_2 where there exists a point $\mathbf{x}_i \rightarrow \mathbf{w}^T \phi(\mathbf{x}_i) > 0$ such that \mathbf{x}_i is in C_1 and there exists a different point $\mathbf{x}_j \rightarrow \mathbf{w}^T \phi(\mathbf{x}_j) < 0$ such that \mathbf{x}_j is in C_2 .

Intuitively, we can see that the decision boundary for the weight vector \mathbf{w} occurs at $\mathbf{w}^T \phi(\mathbf{x}) = 0$

For any point \mathbf{x}_i in C_1 , $|w| \rightarrow \infty$ causes

$$p(C_1 | \phi(\mathbf{x}_i)) = \sigma(\mathbf{w}^T \phi(\mathbf{x}_i)) \rightarrow 1$$

This is because $\sigma(\mathbf{w}^T \phi(\mathbf{x}_i)) = \frac{1}{1 + e^{-\mathbf{w}^T \phi(\mathbf{x}_i)}}$, where large weights will cause the exponent in the denominator to become very negative, rendering that term effectively zero and leaving $\sigma(\mathbf{w}^T \phi(\mathbf{x}_i)) = \frac{1}{1}$.

For any point \mathbf{x}_j in C_2 , $|w| \rightarrow \infty$ causes

$$p(C_2 | \phi(\mathbf{x}_j)) \rightarrow 1$$

This is because

$$\begin{aligned} p(C_2 | \phi(\mathbf{x}_j)) &= 1 - p(C_1 | \phi(\mathbf{x}_j)) \\ &= 1 - \sigma(\mathbf{w}^T \phi(\mathbf{x}_j)) \\ &= 1 - \frac{1}{1 + e^{-\mathbf{w}^T \phi(\mathbf{x}_j)}} \end{aligned}$$

But $\mathbf{w}^T \phi(\mathbf{x}_j)$ is negative (C_2 is represented by predictions less than 0) and has a large magnitude, so we get

$$\begin{aligned} p(C_2 | \phi(\mathbf{x}_j)) &= 1 - \frac{1}{1 + \infty} \\ &= 1 \end{aligned}$$

We therefore see that $|w| \rightarrow \infty$ maximizes the probability of correct classification for both classes – a reason that hard-boundary error functions are more prone to overfitting than soft-boundary ones.

13: PRML Pg. 223, Ex. 4.17

Answer:

Note that $y_k = \frac{e^{a_k}}{\sum_j e^{a_j}}$

For the case $k \neq j$...

$$\begin{aligned} \frac{\partial y_k}{\partial a_j} &= \frac{\partial}{\partial a_j} \frac{e^{a_k}}{\sum_j e^{a_j}} \\ &= e^{a_k} * (-1) * \left(\sum_j e^{a_j} \right)^{-2} * e^{a_j} \\ &= -\frac{e^{a_k}}{\sum_j e^{a_j}} * \frac{e^{a_j}}{\sum_j e^{a_j}} \\ &= -y_k y_j \end{aligned}$$

More generally, we must use the product rule to account for a_j dependency in both the numerator and denominator...

$$\begin{aligned} \frac{\partial y_k}{\partial a_j} &= \frac{e^{a_k}}{\sum_j e^{a_j}} + e^{a_k} * (-1) * \left(\sum_j e^{a_j} \right)^{-2} * e^{a_j} \\ &= \frac{e^{a_k}}{\sum_j e^{a_j}} - \left(\frac{e^{a_k}}{\sum_j e^{a_j}} \right)^2 \\ &= y_k (1 - y_k) \end{aligned}$$

But also note that we may have a vector of outputs \mathbf{y}_k , so in general..

$$\frac{\partial \mathbf{y}_k}{\partial a_j} = \mathbf{y}_k (I_{kj} - y_k)$$

where I_{kj} is the identity matrix.

14: PRML Pg. 223, Ex. 4.18

Answer: Distribute the gradient into the error function sums, noting that the target values do not depend on the weights

$$\nabla_{w_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = - \sum_{n=1}^N \sum_{k=1}^K t_{nk} \nabla_{w_j} \ln y_{nk}$$

Using the chain rule...

$$\nabla_{w_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = - \sum_{n=1}^N \sum_{k=1}^K t_{nk} \frac{1}{y_{nk}} \nabla_{w_j} y_{nk}$$

From PRML Ex. 4.17, and using the chain rule, we can express this as follows:

$$\nabla_{w_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = - \sum_{n=1}^N \sum_{k=1}^K t_{nk} \frac{1}{y_{nk}} y_{nk} (I_{kj} - y_{nj}) \phi_n$$

$$\begin{aligned}
&= - \sum_{n=1}^N \sum_{k=1}^K t_{nk} (I_{kj} - y_{nj}) \phi_n \\
&= \sum_{n=1}^N \sum_{k=1}^K t_{nk} y_{nj} \phi_n - \sum_{n=1}^N \sum_{k=1}^K t_{nk} I_{kj} \phi_n
\end{aligned}$$

We can eliminate the k -dependence in the second-term by acknowledging that anything times the identity matrix is itself projected onto the columns of the identity matrix.

$$\nabla_{w_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = \sum_{n=1}^N \left(\sum_{k=1}^K t_{nk} \right) y_{nj} \phi_n - \sum_{n=1}^N t_{nj} \phi_n$$

For a classification problem, the sum of probabilities of classes for any target example is 1.

$$\begin{aligned}
\nabla_{w_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) &= \sum_{n=1}^N y_{nj} \phi_n - \sum_{n=1}^N t_{nj} \phi_n \\
&= \sum_{n=1}^N (y_{nj} - t_{nj}) \phi_n
\end{aligned}$$