

Draft

7. maj 2014

1 The derivative of prediction or Sensitivity

We wish to find the effect that a datapoint's class has on the predicted class for that datapoint.

$$\frac{\delta \hat{Y}_n}{\delta Y_n} \quad (1.1)$$

Our prediction is

$$\hat{Y}_n = p(y|\bar{x}, \bar{w}) \quad (1.2)$$

where \bar{w} is subject to

$$\frac{\delta L}{\delta \bar{w}} = 0 \quad (1.3)$$

Which means that we have found a locally optimal solution.

We now assume that when we move y by a small amount δy then 1.3 still holds.

Essentially assuming some smoothness around the optimum.

Using this and the fact that 1.3 depends both directly and indirectly on y we see that

$$\begin{aligned} \frac{\delta}{\delta y} \frac{\delta L}{\delta \bar{w}} &= 0 \\ \Downarrow \\ \frac{\delta^2 L}{\delta y \delta \bar{w}} + \frac{\delta^2 L}{\delta \bar{w} \delta \bar{w}^T} \frac{\delta \bar{w}}{\delta y} &= 0 \end{aligned}$$

and from this we can isolate

$$\frac{\delta \bar{w}}{\delta y} = - \left[\frac{\delta^2 L}{\delta \bar{w} \delta \bar{w}^T} \right]^{-1} \frac{\delta^2 L}{\delta y \delta \bar{w}} \quad (1.4)$$

Rewriting 1.1 we get

$$\frac{\delta \hat{Y}_n}{\delta Y_n} = \frac{\delta p(y|\bar{x}, \bar{w})}{\delta Y_n} = \frac{\delta p(y|\bar{x}, \bar{w})}{\delta \bar{w}^T} \frac{\delta \bar{w}}{\delta y} \quad (1.5)$$

And inserting 1.4

$$\frac{\delta p(y|\bar{x}, \bar{w})}{\delta \bar{w}^T} \frac{\delta \bar{w}}{\delta y} = - \frac{\delta p(y|\bar{x}, \bar{w})}{\delta \bar{w}^T} \left[\frac{\delta^2 L}{\delta \bar{w} \delta \bar{w}^T} \right]^{-1} \frac{\delta^2 L}{\delta y \delta \bar{w}} \quad (1.6)$$

2 Randomised algorithm

Uncertainty based on asymptotic likelihood and \bar{w} -distribution

Let \mathcal{L}_∞ be the log-likelihood function for a distribution, now let \mathcal{L}_N denote the log-likelihood function based on N observations from this distribution. Furthermore, let N be a large number, for which $L_N \approx L_\infty$.

$$\mathcal{L}_N = \frac{1}{N} \sum_{n=1}^N \ell_n \quad \bar{w} \text{ s.t. } \frac{\delta \mathcal{L}}{\delta \bar{w}} = \bar{0} \quad (2.1)$$

Where ℓ_n is the log-likelihood of the n^{th} observation. And \bar{w} is the optimal (**true?**) weights for the distribution, then we combine the expressions from (2.1), such that for the optimal weights the following must be fulfilled:

$$\frac{1}{N} \sum_{n=1}^N \frac{\delta \ell_n}{\delta \bar{w}} = 0 \quad (2.2)$$

(Skal vi lige skrive lidt om at $\Delta w = w - w_0$ og er en lille forskydelse i vægtene? Eller er det en lille forskydelse?) For each of the N observations, we can write the log-likelihood of the n^{th} observation as:

$$\ell_n(\Delta \bar{w}) = \ell_n(\bar{w}_0) + \left. \frac{\delta \ell_n}{\delta \bar{w}} \right|_{\bar{w}_0} \Delta \bar{w} + \frac{1}{2} \text{Tr} \left[\left. \frac{\delta^2 \ell_n}{\delta \bar{w} \delta \bar{w}^T} \right|_{\bar{w}_0} \Delta \bar{w} \Delta \bar{w}^T \right] \quad (2.3)$$

Or for the entire log-likelihood function:

$$\mathcal{L}_N(\Delta \bar{w}) = \mathcal{L}_\infty(\bar{w}_0) + \left(\left. \frac{\delta \mathcal{L}_N}{\delta \bar{w}} \right|_{\bar{w}_0} \right)^T \cdot \Delta \bar{w} + \frac{1}{2} \Delta \bar{w}^T \left(\left. \frac{\delta^2 \mathcal{L}_N}{\delta \bar{w} \delta \bar{w}^T} \right|_{\bar{w}_0} \right) \Delta \bar{w} + R \quad (2.4)$$

Where R is the error of the approximation. Furthermore, we define the functions $\bar{\bar{H}}_N = \left. \frac{\delta^2 \mathcal{L}_N}{\delta \bar{w} \delta \bar{w}^T} \right|_{\bar{w}_0}$, and $\bar{g} = \left. \frac{\delta \mathcal{L}_N}{\delta \bar{w}} \right|_{\bar{w}_0}$. And evaluate the condition on \bar{w} , stated in (2.1):

$$\frac{\delta \mathcal{L}_N}{\delta \bar{w}} = \bar{g}_N + \bar{\bar{H}}_N \Delta \bar{w} = \bar{0} \quad (2.5)$$

We replace $\Delta \bar{w}$ with $\hat{\Delta \bar{w}}$ as N is a finite number, thus only approximating $\Delta \bar{w}$. Isolating $\hat{\Delta \bar{w}}$, and using Ljung [REFERENCE?], we get:

$$\hat{\Delta \bar{w}} = -\bar{\bar{H}}_N^{-1} \cdot \bar{g}_N \stackrel{\text{Ljung}}{=} -\bar{\bar{H}}_0^{-1} \cdot \bar{g}_{\bar{w}}(\bar{w}_0) \quad (2.6)$$

2.1 Covariance of \bar{w} - distribution

(Forklaring af at H_0 er uafhængig af datasæt, mens g nu er afhængig af w evalueret i w_0) Besides getting an estimate for $\Delta\hat{w}$, we can find the mean of the distribution:

$$\langle \Delta\hat{w} \rangle = -\bar{H}_0^{-1} \langle \bar{g}_w \rangle (\bar{w}_0) = 0$$

As $\delta\bar{w} = \bar{w} - \bar{w}_0$??mistet tråden?

2.1 Covariance of \bar{w} - distribution

$$\langle \delta\bar{w}\delta\bar{w}^T \rangle_{D_N} = \left\langle \bar{H}^{-1} \bar{g} \bar{g}^T \bar{H}^{-1} \right\rangle \stackrel{Ljung}{=} \bar{H}_0^{-1} \langle \bar{g} \bar{g}^T \rangle \bar{H}_0^{-1} + R' \quad (2.7)$$

With error $R' = O(\frac{1}{N}) \approx 0$, for large N . We look at the covariance of the gradient function

$$\begin{aligned} \langle \bar{g} \bar{g}^T \rangle_N &= \frac{1}{N^2} \sum_{n,n'=1}^N \left\langle \frac{\delta \ell_n}{\delta \bar{w}} \bigg|_{\bar{w}_0} \frac{\delta \ell_{n'}}{\delta \bar{w}} \bigg|_{\bar{w}_0} \right\rangle \\ &= \frac{1}{N^2} \left(\sum_{n \neq n'} \left\langle \frac{\delta \ell_n}{\delta \bar{w}} \bigg|_{\bar{w}_0} \right\rangle \cdot \left\langle \frac{\delta \ell_{n'}}{\delta \bar{w}} \bigg|_{\bar{w}_0} \right\rangle + \sum_{n=1}^N \left\langle \frac{\delta \ell_n}{\delta \bar{w}} \bigg|_{\bar{w}_0} \frac{\delta \ell_n}{\delta \bar{w}^T} \bigg|_{\bar{w}_0} \right\rangle \right) \end{aligned} \quad (2.8)$$

$$(2.9)$$

Due to the assumption of independence, only the N diagonal elements are non-zero. So;

$$\langle \bar{g} \bar{g}^T \rangle_N = \frac{1}{N} \left\langle \frac{\delta \mathcal{L}}{\delta \bar{w}} \bigg|_{\bar{w}_0} \frac{\delta \mathcal{L}}{\delta \bar{w}^T} \bigg|_{\bar{w}_0} \right\rangle \quad (2.10)$$

2.2 Proof that $\left\langle \frac{\delta \mathcal{L}}{\delta \bar{w}} \bigg|_{\bar{w}_0} \frac{\delta \mathcal{L}}{\delta \bar{w}^T} \bigg|_{\bar{w}_0} \right\rangle = \bar{H}_0$

Tekst test

$$\langle \bar{g} \bar{g}^T \rangle_N = \frac{1}{N^2} \sum_{n=1}^N \int_{\Omega} \frac{\delta \ell_n(\bar{x})}{\delta \bar{w}} \bigg|_{\bar{w}_0} \frac{\delta \ell_n(\bar{x})}{\delta \bar{w}^T} \bigg|_{\bar{w}_0} p(\bar{x}) \delta x \quad (2.11)$$

From (2.10) and (2.11), and setting $\ell_n(\bar{x}) = p(\bar{x})$:

$$\bar{H} \bigg|_{\bar{w}_0} = \frac{1}{N} \sum_{n=1}^N \int_{\Omega} \frac{\delta}{\delta \bar{w} \delta \bar{w}^T} - \log p(\bar{x}|\bar{w}) p(\bar{x}) \delta x \quad (2.12)$$

$$= \frac{1}{N} \sum_{n=1}^N \int_{\Omega} -\frac{\delta}{\delta \bar{w}} \frac{1}{p(\bar{x})} \frac{\delta}{\delta \bar{w}^T} p(\bar{x}|\bar{w}) p(\bar{x}) \delta \bar{x} \quad (2.13)$$

$$(2.14)$$

Now if $p(\bar{x}|\bar{w}_0) = p(x)$, then