

# fake title

Jesper Løve Hinrich  
912933169

6. 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  $\delta y$  then ?? still holds.

Essentially assuming some smoothness around the optimum.

Using this and the fact that ?? 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 ?? 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 ??

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

## 2 Randomised algorithm

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

Let  $\mathcal{L}_\infty$  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  $\ell_n$  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 (??), 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:

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

Or for the entire log-likelihood function:

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

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

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

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:

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

(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:

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

**Nu skal vi prøve at kombinere flere datasets**