

## Vedlegg B

### I. ADAM

Adam (?) er en optimeringsalgoritme som er en tilpasning av SGD hvor læringsraten  $\eta$  varierer gjennom iterasjonene. I Adam inkluderes gradienten fra forrige iterasjon i første og andre potens i tillegg, på en måte som kan betraktes som en tilpasning av læringsraten.

Adam har flere metaparametere,  $\eta$ ,  $\rho_1$ ,  $\rho_2$ , og vi ser at denne metoden benytter både gradienten i første potens ( $\mathbf{m}$ ), med minne fra forrige iterasjon, og gradienten i andre potens ( $\mathbf{s}$ ) for å justere steget i  $\mathbf{w}$ :

$$\mathbf{m}^{(i)} = \rho_1 \mathbf{m}^{(i-1)} + (1 - \rho_1) \mathbf{g}^{(i)}$$

$$\mathbf{s}^{(i)} = \rho_2 \mathbf{s}^{(i-1)} + (1 - \rho_2) (\mathbf{g}^{(i)})^2$$

$$\mathbf{m}^{(i)} = \frac{\mathbf{m}^{(i)}}{1 - \rho_1^i}$$

$$\mathbf{s}^{(i)} = \frac{\mathbf{s}^{(i)}}{1 - \rho_2^i}$$

$$\mathbf{w}^{(i+1)} = \mathbf{w}^{(i)} - \eta^{(i)} \frac{\mathbf{m}^{(i)}}{\sqrt{\mathbf{s}^{(i)} + \delta}}.$$

$\delta$  er her et lite tall som legges til for å unngå problemer med å dele på 0. Og igjen, vil det være et helt tilsvarende sett med ligninger for  $\mathbf{b}$ , og det er vanlig å benytte de samme metaparameterne for å optimere både  $\mathbf{w}$  og  $\mathbf{b}$ .