

Vedlegg A

I. BESKRIVELSE AV LOGSITISK REGRESJON

A. Binær logistisk regresjon

For å modellere sannsynlighet med binær logistisk regresjon, brukes sigmoidfunksjonen som tar verdier mellom 0 og 1:

$$\sigma(z) = \frac{1}{1 + \exp(-z)} = \frac{\exp z}{1 + \exp z}.$$

Her er z representert som $z = \beta_0 + \beta_1 x$, der β_0 og β_1 er parametere som må tilpasses i modellen. Da blir

$$p(x) = \frac{\exp(\beta_0 + \beta_1 x)}{1 + \exp(\beta_0 + \beta_1 x)}$$

som kan generaliseres til

$$p(x) = \frac{\exp(\mathbf{w}x + \mathbf{b})}{1 + \exp(\mathbf{w}x + \mathbf{b})},$$

i tilfellet hvor inndataene våre består av mer enn én parameter. I et datasett bestående av sammenhørende verdier \mathbf{x}_i og y_i , vil da sannsynligheten for gitt utfall være

$$P = \prod_{i=1}^n [p(y_i = 1 | \mathbf{x}_i, \mathbf{w}, \mathbf{b})]^{y_i} [1 - p(y_i = 1 | \mathbf{x}_i, \mathbf{w}, \mathbf{b})]^{1-y_i}.$$

Vi ønsker å finne verdiene \mathbf{w} og \mathbf{b} som gir størst mulig sannsynlighet, dvs at modellen vår passer best mulig med det som er observert. For å gjøre det enklere, så tar vi den naturlige logaritmen av denne funksjonen

$$\begin{aligned} \log(P) &= \sum_{i=1}^n (y_i \log p(y_i = 1 | \mathbf{x}_i, \mathbf{w}, \mathbf{b}) \\ &\quad + (1 - y_i) \log [1 - p(y_i = 1 | \mathbf{x}_i, \mathbf{w}, \mathbf{b})]). \end{aligned}$$

Hvis vi så setter inn $p(x)$ i denne får vi et uttrykk som kalles for *cross entropy* som benyttes som kost-funksjon i logistisk regresjon:

$$\mathcal{C}(\mathbf{w}, \mathbf{b}) = - \sum_{i=1}^n (y_i (\mathbf{w}x_i + \mathbf{b}) - \log(1 + \exp(\mathbf{w}x_i + \mathbf{b}))).$$

Optimering, eller trening, av modellen består så i å finne det beste settet med parametere \mathbf{w} og \mathbf{b} . Deretter kan modellen benyttes til å klassifisere ukjente tilfeller gjennom å beregne sannsynligheten $p(x)$ og en på forhånd bestemt terskelverdi for å angi hvilken klasse den hører hjemme i.

B. Klassifisering med mer enn to klasser

Dersom klassifiseringsproblemet har mer enn to klasser, må modellen over tilpasses for å ta høyde for mer enn to utfall. Sannsynligheten for et gitt utfall blir nå

$$P = \prod_{i=1}^n \prod_{k=1}^K p_k(\mathbf{x}_i, \mathbf{w}, \mathbf{b})^{y_i},$$

hvor det er totalt K klasser. Sannsynlighetsmodellen må også tilpasses at det er flere klasser og dette gjøres gjennom å erstatte sigmoidfunksjonen med *softmax*:

$$p_k(\mathbf{x}_i) = \frac{\exp(\mathbf{w}_k x + \mathbf{b}_k)}{\sum_{l=1}^{K-1} \exp(\mathbf{w}_l x + \mathbf{b}_l)}.$$

Totalt sett får vi da til slutt en kostfunksjon gitt ved:

$$\mathcal{C}(\mathbf{w}, \mathbf{b}) = - \sum_{i=1}^n \sum_{k=1}^K \log \frac{\exp(\mathbf{w}_k x + \mathbf{b}_k)}{\sum_{l=1}^{K-1} \exp(\mathbf{w}_l x + \mathbf{b}_l)}.$$

Denne kostfunksjonen kan så benyttes for å finne optimale verdier av modellparameterne. Men for å gjøre prediskjoner i flerklasselklassifisering, klassifisering med mer enn to utfall, må det beregnes en sannsynlighet for hver av klassene. Deretter kan man f.eks. benytte såkalt *one hot encoding* som vil si at man plasserer tilfellet i den klassen som har høyest sannsynlighet. På denne måten vil man i alle tilfeller tildele en klasse. Modellen vi da predikere kun én klasse selv om det var flere klasser som hadde høy sannsynlighet og også selv om ingen av klassene hadde høy sannsynlighet.