# Artificial neural networks

- Feedforward NN with one hidden layer:

$$\hat{\mathbf{y}} = \varphi(\mathbf{W}^{(2)}\sigma^{(1)}(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)}) = \varphi\left( \sum_{k=1}^{K} W_{kl}^{(2)}\, \sigma^{(1)}\left( \sum_{j=1}^{J} W_{jk}^{(1)}x_j + b_j^{(1)} \right) + b_k^{(2)} \right)$$

$\mathbf{x}$ = input layer

$\hat{\mathbf{y}}$ = output prediction layer

$\boldsymbol{\theta}$ = parameters to estimate = $\{\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{W}^{(2)}, \mathbf{b}^{(2)}\}$

$$\sigma(\mathbf{z}) = \begin{cases} \max(\mathbf{0}, \mathbf{z}) & \text{relu, defacto standard} \\ \left(1 + e^{-\mathbf{z}}\right)^{-1} & \text{sigmoid, old school} \\ \text{many} & \text{variations on these and others} \end{cases}$$

$$\varphi(\mathbf{z}) = \begin{cases} h\tan \mathbf{z} & regression \\ \dfrac{e^{\mathbf{z}}}{\sum_{\mathbf{z}} e^{\mathbf{z}}} & classification \end{cases}$$

# NN parameter estimation for regression

- Model:
$$\mathbf{Y} = \hat{\mathbf{Y}} + \boldsymbol{\epsilon} = f(\mathbf{X}; \boldsymbol{\theta}) + \boldsymbol{\epsilon} \qquad where \qquad f(\mathbf{X}; \boldsymbol{\theta}) \text{ expresses our neural network}$$

$$\boldsymbol{\epsilon} \sim N(\mathbf{Y} - f(\mathbf{X}; \boldsymbol{\theta}), \boldsymbol{\Sigma})$$

$$= \sqrt{\frac{1}{(2\pi)^N \det \boldsymbol{\Sigma}}} e^{-\frac{1}{2}(\mathbf{X} - f(\mathbf{X}; \boldsymbol{\theta}))^T \boldsymbol{\Sigma}^{-1}(\mathbf{X} - f(\mathbf{X}; \boldsymbol{\theta}))}$$

- Optimization:
$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \; -\log P(\mathbf{Y} \,|\, \mathbf{X}; \boldsymbol{\theta}) \qquad := P(\mathbf{Y} \,|\, \mathbf{X}; \boldsymbol{\theta}) \; \leftarrow \; \mathscr{L}(\mathbf{D}; \boldsymbol{\theta})$$

$$= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \; -\log\left(\sqrt{\frac{1}{(2\pi)^N \det \boldsymbol{\Sigma}}} \exp\left[-\frac{1}{2}\big(\mathbf{X} - f(\mathbf{X}; \boldsymbol{\theta})\big)^T \boldsymbol{\Sigma}^{-1}\big(\mathbf{X} - f(\mathbf{X}; \boldsymbol{\theta})\big)\right]\right)$$

$$= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \; -\frac{1}{2}\log\big((2\pi)^N \det \boldsymbol{\Sigma}\big) - \frac{1}{2}\big(\mathbf{X} - f(\mathbf{X}; \boldsymbol{\theta})\big)^T \boldsymbol{\Sigma}^{-1}\big(\mathbf{X} - f(\mathbf{X}; \boldsymbol{\theta})\big)$$

$$= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \; -\big(\mathbf{X} - f(\mathbf{X}; \boldsymbol{\theta})\big)^T\big(\mathbf{X} - f(\mathbf{X}; \boldsymbol{\theta})\big)$$

$$= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \; -\sum_{i=1}^{M}\big(\mathbf{y}^{(i)} - f(\mathbf{x}^{(i)}; \boldsymbol{\theta})\big)^T\big(\mathbf{y}^{(i)} - f(\mathbf{x}^{(i)}; \boldsymbol{\theta})\big)$$

$$= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \; -\sum_{i=1}^{M}\sum_{j=1}^{N}\big(\mathbf{y}^{(i)}_j - f(\mathbf{x}^{(i)}; \boldsymbol{\theta})_j\big)^2 \qquad \leftarrow \text{ least squares}$$

$\boldsymbol{\Sigma}$ is diagonal, strictly positive, independent of $\mathbf{x}$ ; it doesn't affect $\hat{\boldsymbol{\theta}}$