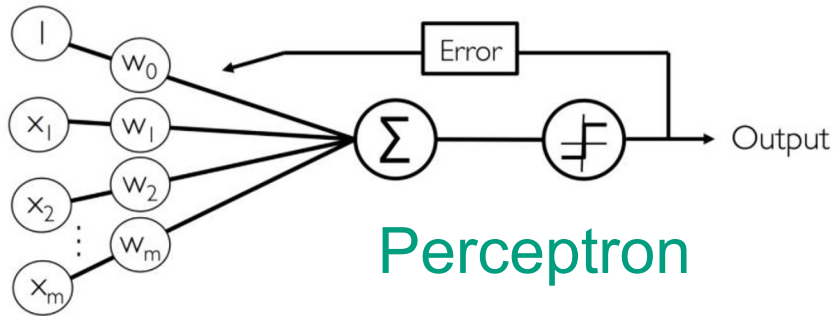


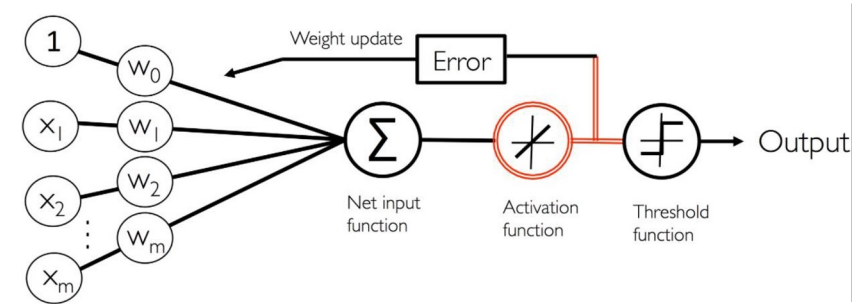
Quadratic loss/cost function

Recap

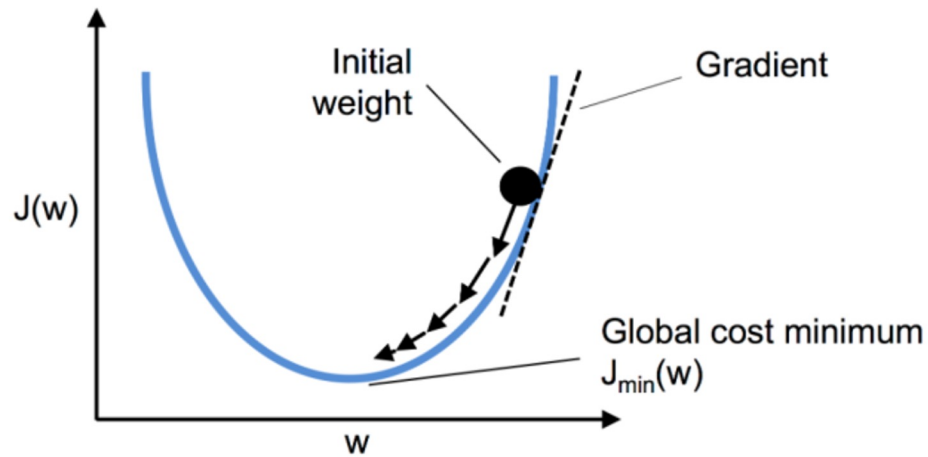
$$J(\mathbf{w}) = \frac{1}{2} \sum_i \left(y^{(i)} - \phi(z^{(i)}) \right)^2$$



Perceptron



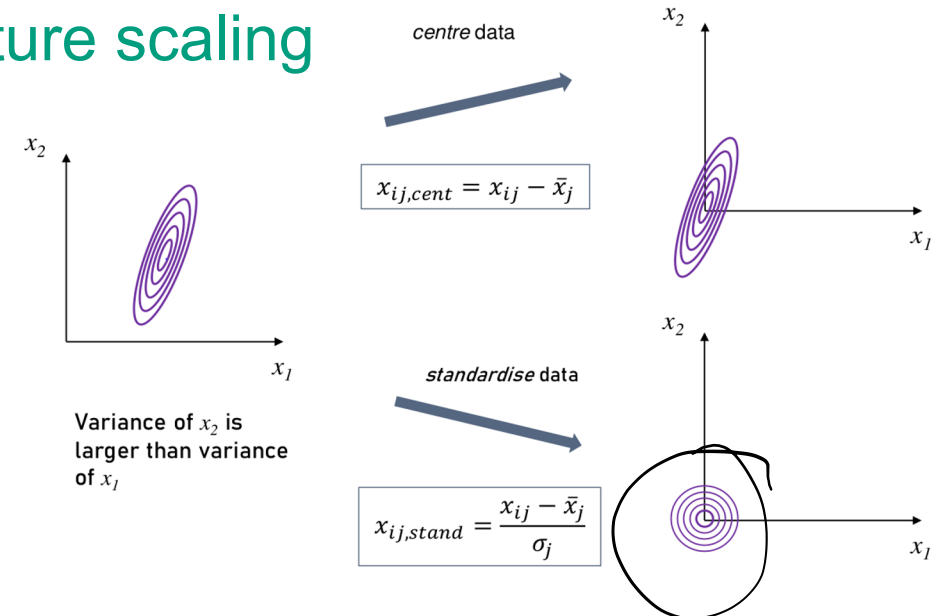
Adaptive linear neuron (Adaline)

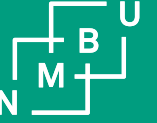


$$\Delta \mathbf{w} = -\eta \nabla J(\mathbf{w})$$

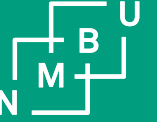
Gradient descent

Feature scaling





Linear regression & logistic regression



Basic notation (repetition)



Vectors and Matrices

We often represent raw (numeric) data as vectors and matrices

Example: Iris data can be represented as a 150 by 4 matrix: $\mathbf{X} \in \mathbb{R}^{150 \times 4}$

- Superscript means i-th training sample
- Subscript means j-th feature (dimension)
- Lowercase boldface \rightarrow vectors ($\mathbf{x} \in \mathbb{R}^{n \times 1}$)
- Uppercase boldface \rightarrow matrices ($\mathbf{X} \in \mathbb{R}^{m \times n}$)
- Single element in a vector $x^{(i)}$
- Single element in a matrix $x_j^{(i)}$

$$\mathbf{X} = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & x_3^{(1)} & x_4^{(1)} \\ x_1^{(2)} & x_2^{(2)} & x_3^{(2)} & x_4^{(2)} \\ \vdots & \vdots & \vdots & \vdots \\ x_1^{(150)} & x_2^{(150)} & x_3^{(150)} & x_4^{(150)} \end{bmatrix}$$



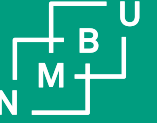
Vectors and Matrices

Row vectors (e.g. one row in \mathbf{X} , “flower” sample in Iris data set)

$$\mathbf{x}^{(i)} = \begin{bmatrix} x_1^{(i)} & x_2^{(i)} & x_3^{(i)} & x_4^{(i)} \end{bmatrix} \quad \mathbf{x}^i \in \mathbb{R}^{1 \times 4}$$

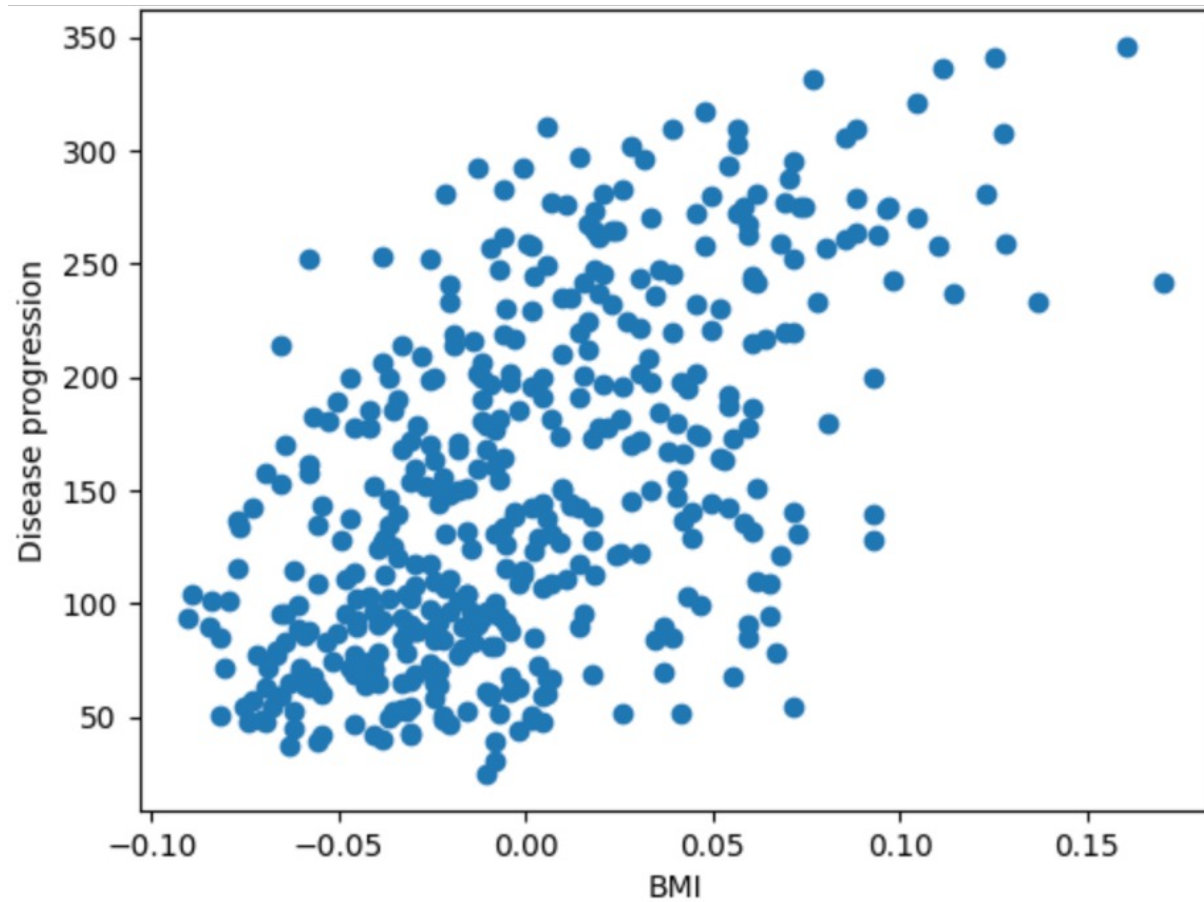
Column vectors (e.g. one column in \mathbf{X} , one feature)

$$\mathbf{x}_j = \begin{bmatrix} x_j^{(1)} \\ x_j^{(2)} \\ \vdots \\ x_j^{(150)} \end{bmatrix} \quad \mathbf{x}_j \in \mathbb{R}^{150 \times 1}$$

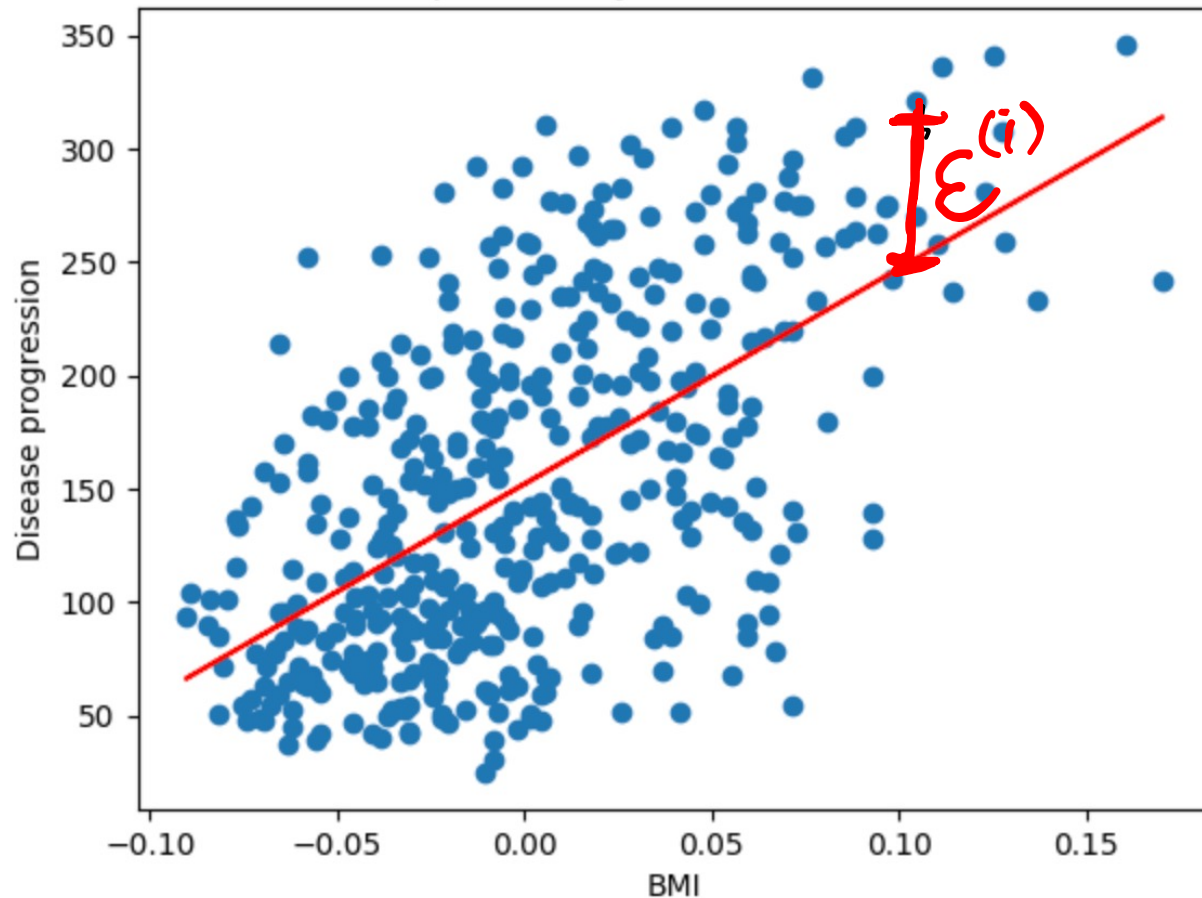


Simple linear regression (and least squares)

Linear regression



Linear regression



Model (parametric)

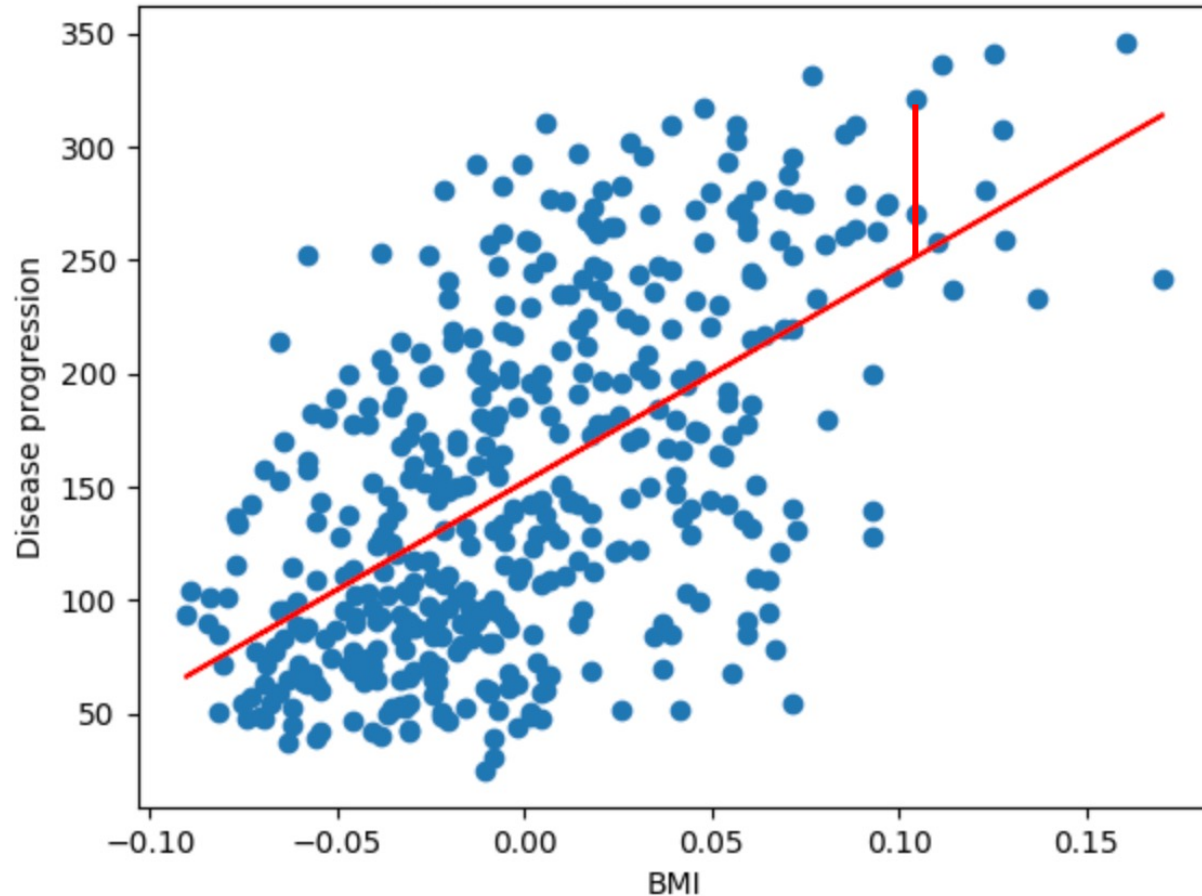
$$f(x) = \theta_0 + \theta_1 x \quad y = X\theta$$

Error

$$\epsilon^{(i)} = y^{(i)} - (\theta_0 + \theta_1 x^{(i)})$$

Linear regression

$$\varepsilon^{(i)} = y^{(i)} - (\theta_0 + \theta_1 x^{(i)})$$



Generalization to m features

$$\varepsilon^{(i)} = y^{(i)} - \mathbf{x}^{(i)} \boldsymbol{\theta}$$

column vector

$$\mathbf{x}^{(i)} = [1, x_1] \quad \boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}$$

row-vector

$$\mathbf{x}^{(i)} \boldsymbol{\theta} := \sum_{j=0}^n \varepsilon_j \theta_j$$

$\theta_0 + \theta_1 x_1$

Alternative equivalent parameter representation

$$\mathbf{x}^{(i)} \boldsymbol{\theta} = b + \mathbf{x}^{(i)} \mathbf{w}$$

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_m \end{bmatrix} := \begin{bmatrix} b \\ \mathbf{w} \end{bmatrix}$$

$$\mathbf{w} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_m \end{bmatrix}, \quad b = \theta_0$$

“weights”

“normal”

“bias”

“offset”

(hyperplanes)



Linear regression

Error in matrix notation

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_m \end{bmatrix}$$

$$\varepsilon = \mathbf{y} - \mathbf{X}\theta$$

u samples

$$\varepsilon = \begin{bmatrix} \varepsilon^{(1)} \\ \varepsilon^{(2)} \\ \varepsilon^{(3)} \\ \vdots \\ \varepsilon^{(n)} \end{bmatrix}$$

$$\mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

$$\mathbf{X} = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(n)} \end{bmatrix} = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_m^{(1)} \\ 1 & x_1^{(2)} & \dots & x_m^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \dots & x_m^{(n)} \end{bmatrix}$$

Linear regression

Goal: Minimize the error

We have one error for each sample, how to measure a global error?

Quadratic / Squared-error loss function

$$L = \underline{\boldsymbol{\epsilon}^T \boldsymbol{\epsilon}} = ||\boldsymbol{\epsilon}||_2^2 = \sum_{i=1}^n \epsilon^{(i)2}$$

$$||\boldsymbol{\epsilon}||_2 = \sqrt{\epsilon_1^2 + \dots + \epsilon_n^2}$$

Least squares solution

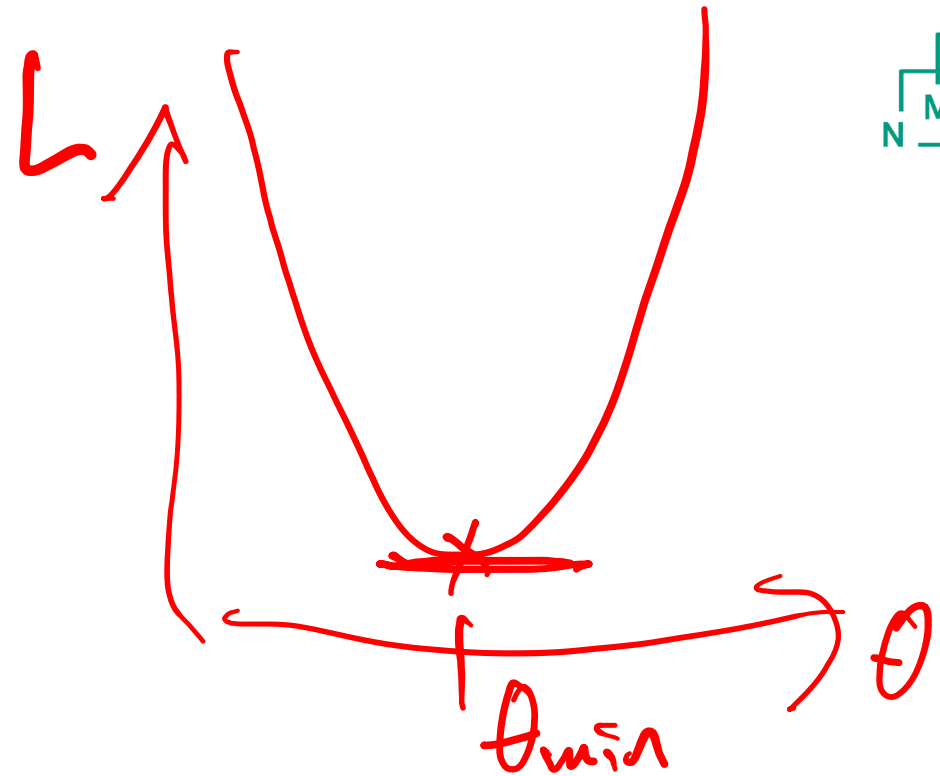
Goal: Minimize the sum of squared errors

$$L = \boldsymbol{\epsilon}^T \boldsymbol{\epsilon} = \|\boldsymbol{\epsilon}\|_2^2 = \sum_{i=1}^n \epsilon^{(i)2}$$

Insert error

$$L(\boldsymbol{\theta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})$$

$$= \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\boldsymbol{\theta} - \boldsymbol{\theta}^T \mathbf{X}^T \mathbf{y} + \boldsymbol{\theta}^T \mathbf{X}^T \mathbf{X}\boldsymbol{\theta}$$



$$\underline{\epsilon} = \mathbf{y} - \mathbf{X}\boldsymbol{\theta}$$

Least squares solution

Goal: Minimize the sum of squared errors

$$L(\boldsymbol{\theta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) = \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\boldsymbol{\theta} - \boldsymbol{\theta}^T \mathbf{X}^T \mathbf{y} + \boldsymbol{\theta}^T \mathbf{X}^T \mathbf{X}\boldsymbol{\theta}$$

$$\frac{\partial L}{\partial \boldsymbol{\theta}} = \frac{\partial \left(\mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\boldsymbol{\theta} - \boldsymbol{\theta}^T \mathbf{X}^T \mathbf{y} + \boldsymbol{\theta}^T \mathbf{X}^T \mathbf{X}\boldsymbol{\theta} \right)}{\partial \boldsymbol{\theta}} = \mathbf{0}$$

$$= -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X}\boldsymbol{\theta} = \mathbf{0}$$

Least squares solution

$$\frac{\partial L}{\partial \theta} = -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \theta = 0$$

$$\mathbf{X}^T \mathbf{X} \theta = \mathbf{X}^T \mathbf{y}$$

(Normal equations)

$$\theta = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Least squares solution



Least squares solution

Learning / Training / Fitting

$$\boldsymbol{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Prediction

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta}$$

Remark: remember that we added a “1” feature to \mathbf{X}

```
lin_regression.ipynb
```


Linear regression and Adaline?

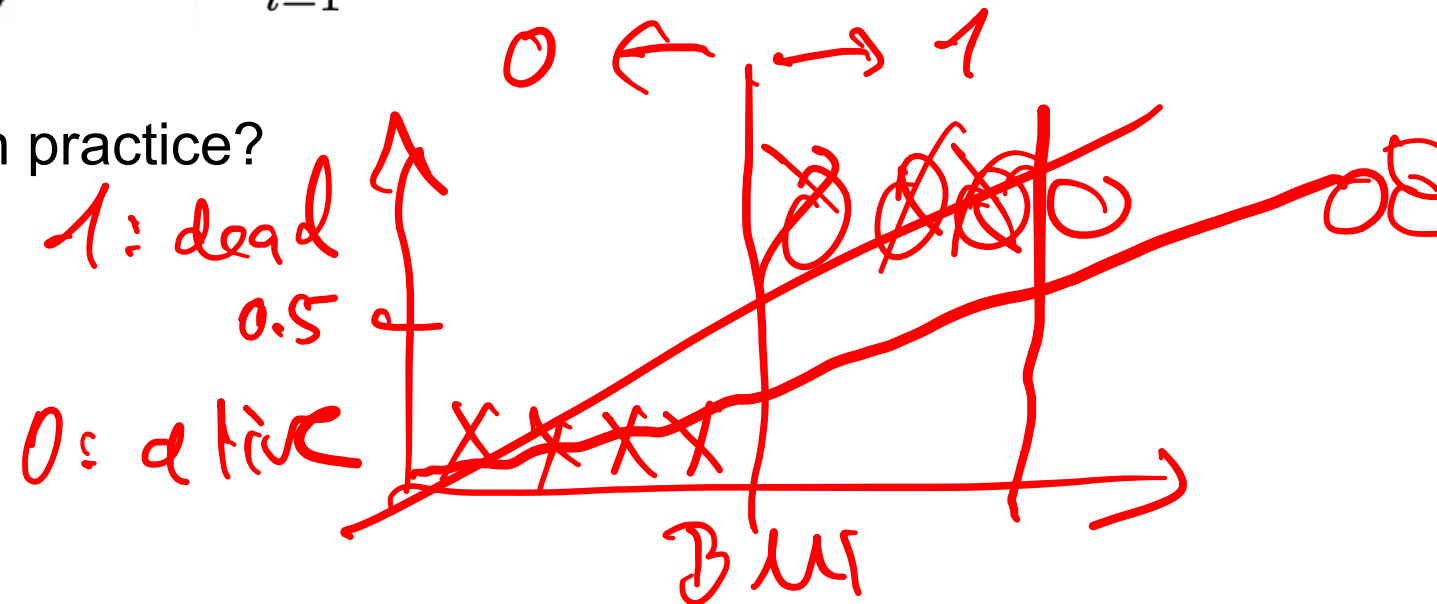
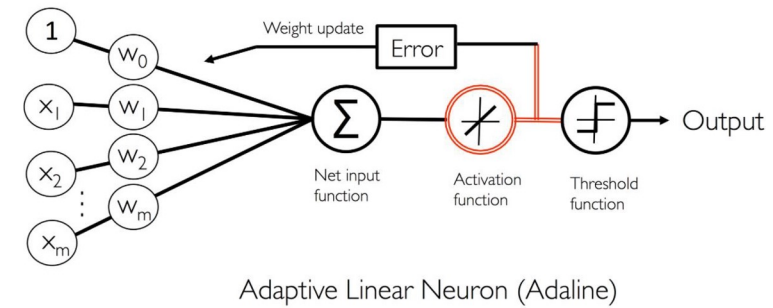
→ Adaline is linear regression with a threshold

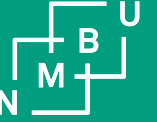
Same loss function:

$$L = \sum_{i=1}^n \varepsilon^{(i)2} = \sum_{i=1}^n \left(y^{(i)} - \sum_{k=0}^m x_k^{(i)} \theta_k \right)^2 = \sum_{i=1}^n \left(y^{(i)} - \mathbf{x}^{(i)} \boldsymbol{\theta} \right)^2$$

Why is Adaline not always so good in practice?

- sensitive to outliers





Linear regression again

Probabilistic interpretation

Probabilistic interpretation

Linear model

$$y^{(i)} = \mathbf{x}^{(i)} \boldsymbol{\theta} + \varepsilon^{(i)}$$

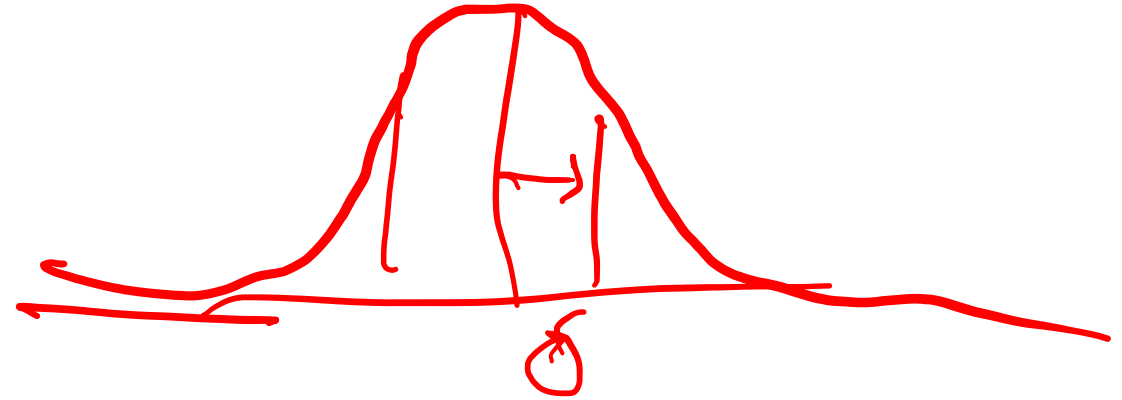
Error model (assumption)

$$\underline{\varepsilon^{(i)}} \sim \underline{\mathcal{N}(0, \sigma^2)}$$

$\varepsilon^{(i)}$ are i.i.d.

$$\underline{p(\varepsilon^{(i)}) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{\varepsilon^{(i)2}}{2\sigma^2}\right)}$$

independent and identically distributed



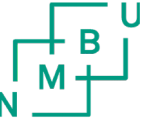
Probabilistic interpretation

Probability of a single outcome, given the sample, and parametrized by θ

$$P(y^{(i)} \mid x^{(i)}; \theta) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y^{(i)} - \mathbf{x}^{(i)}\theta)^2}{2\sigma^2}\right)$$

Likelihood of θ (defined in terms of probability of the data)

$$\mathcal{L}(\theta) = P(\mathbf{y} \mid \mathbf{X}; \theta)$$



Probabilistic interpretation

Goal: We want to maximize the likelihood of our parameters

$$\mathcal{L}(\boldsymbol{\theta}) = P(\mathbf{y} \mid \mathbf{X}; \boldsymbol{\theta})$$

$$= \prod_{i=1}^n P(y^{(i)} \mid x^{(i)}; \boldsymbol{\theta}), \quad (\text{since } \varepsilon^{(i)} \text{ are i.i.d.})$$

$$= \prod_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} \exp \left(-\frac{(y^{(i)} - \mathbf{x}^{(i)} \boldsymbol{\theta})^2}{2\sigma^2} \right)$$

$$\prod_{i=1}^n \alpha_i = \alpha_1 \alpha_2 \alpha_3 \cdots \alpha_n$$

Probabilistic interpretation

$$\mathcal{L}(\boldsymbol{\theta}) = \prod_{i=1}^n P(y^{(i)} | x^{(i)}; \boldsymbol{\theta})$$

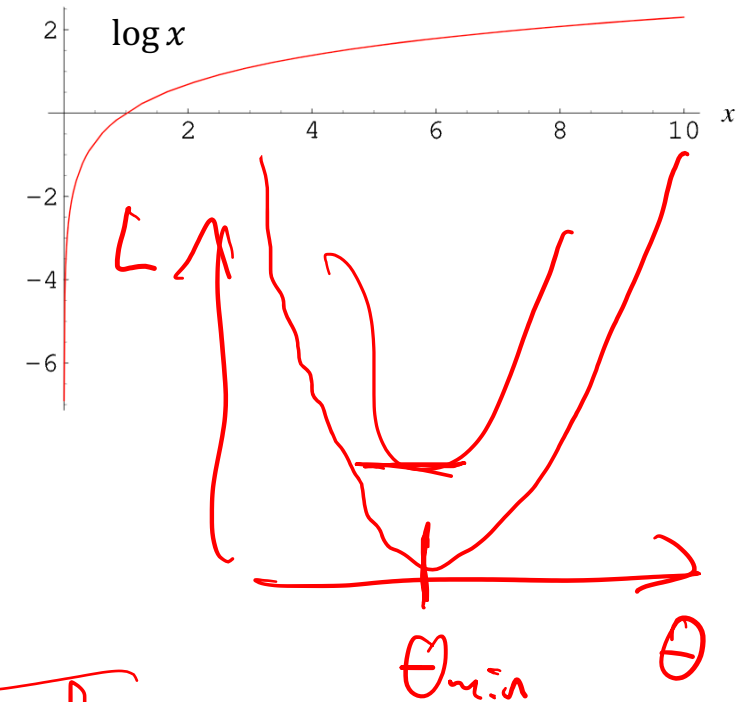
In practice it's often beneficial to look at the log-likelihood

Since the natural logarithm (\log) is a strictly monotone function, likelihood and log-likelihood attain maximum at the same $\boldsymbol{\theta}$

$$\ell(\boldsymbol{\theta}) := \log \mathcal{L}(\boldsymbol{\theta}) = \log \prod_{i=1}^n P(y^{(i)} | x^{(i)}; \boldsymbol{\theta})$$

$$= \log \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} \exp \left(-\frac{(y^{(i)} - \mathbf{x}^{(i)}\boldsymbol{\theta})^2}{2\sigma^2} \right)$$

$$= n \log \frac{1}{\sigma\sqrt{2\pi}} - \frac{1}{2\sigma^2} \sum_{i=1}^n (y^{(i)} - \mathbf{x}^{(i)}\boldsymbol{\theta})^2$$



Probabilistic interpretation

Maximizing a function is the same as minimizing the negative function

$$\ell(\boldsymbol{\theta}) := \log \mathcal{L}(\boldsymbol{\theta}) = n \log \frac{1}{\sigma \sqrt{2\pi}} - \frac{1}{2\sigma^2} \sum_{i=1}^n \left(y^{(i)} - \mathbf{x}^{(i)} \boldsymbol{\theta} \right)^2$$

New goal: Minimize negative log-likelihood

(leave away scaling factors and constant)

$$L(\boldsymbol{\theta}) = \sum_{i=1}^n \left(y^{(i)} - \mathbf{x}^{(i)} \boldsymbol{\theta} \right)^2 = \sum_{i=1}^n \varepsilon^{(i)2}$$

→ For example, solve with least squares

Probabilistic interpretation

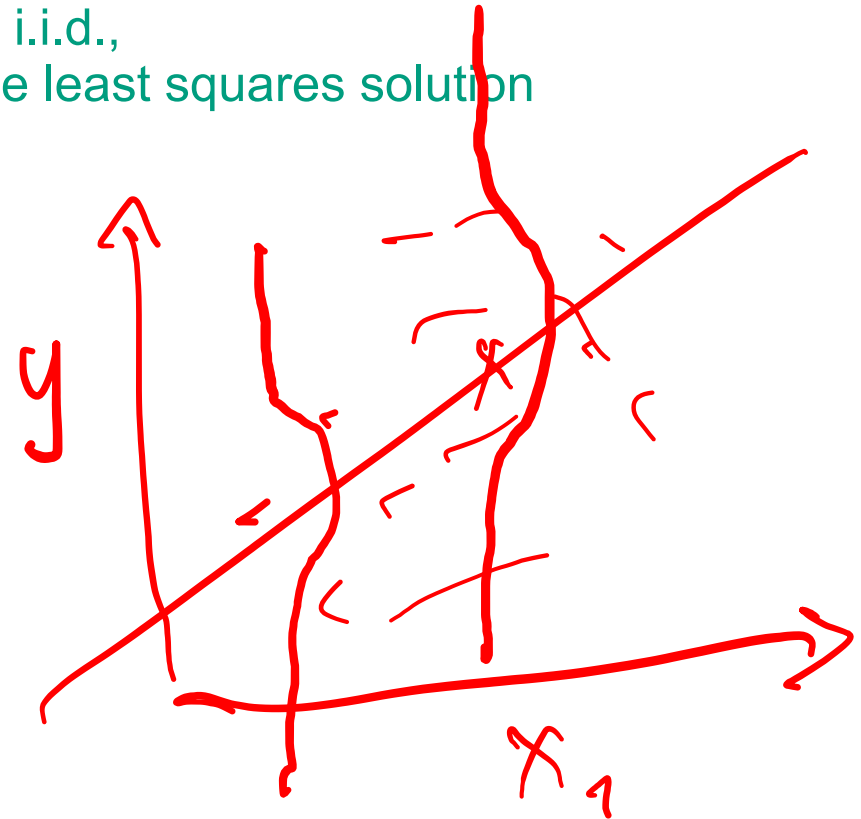
Under the assumption that the errors are Gaussian and i.i.d., the **maximum likelihood estimator** for θ is given by the least squares solution

$$\theta_{\text{MLE}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

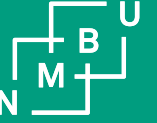
What is the variance σ^2 ?

$$\frac{\partial \ell(\theta, \sigma)}{\partial \sigma} = -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n \left(y^{(i)} - \mathbf{x}^{(i)T} \theta \right)^2 = 0$$

$$\rightarrow \sigma^2 = \frac{1}{n} \sum_{i=1}^n \left(y^{(i)} - \mathbf{x}^{(i)T} \theta \right)^2 = \frac{1}{n} \sum_{i=1}^n \varepsilon^{(i)2}$$



variance is the mean squared error



Logistic regression



Logistic regression, a binary classifier

Labels $y^{(i)} \in \{0, 1\}$

Probability of the data

$$p := P(y^{(i)} = 1 \mid \mathbf{x}^{(i)}) \rightarrow P(y^{(i)} = 0 \mid \mathbf{x}^{(i)}) = 1 - p.$$

What are the odds?

Useful in practice: log-odds

$$\frac{p}{1-p}$$

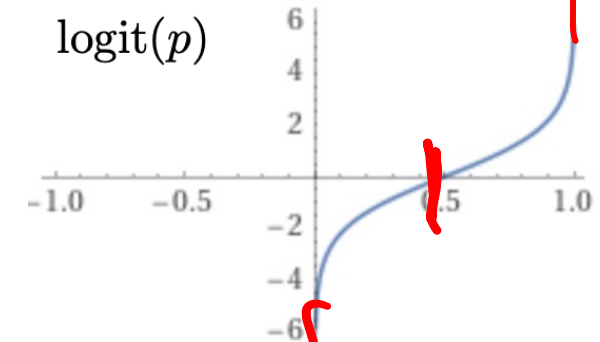
$$\text{logit}(p) := \log \frac{p}{1-p}, \quad \text{logit} : (0, 1) \rightarrow \mathbb{R}$$

Logistic regression, the model

The log-odds are modelled by a linear function

$$\text{logit}(p) = \mathbf{x}^{(i)} \boldsymbol{\theta} = b + \mathbf{x}^{(i)} \mathbf{w}$$

$$\text{logit}(p) := \log \frac{p}{1-p}, \quad \text{logit} : (0, 1) \rightarrow \mathbb{R}$$



log-odds

$$\sigma(\text{logit}(p)) = p \quad \sigma(\mathbf{x}^{(i)} \boldsymbol{\theta}) = p$$

Logistic regression, the model

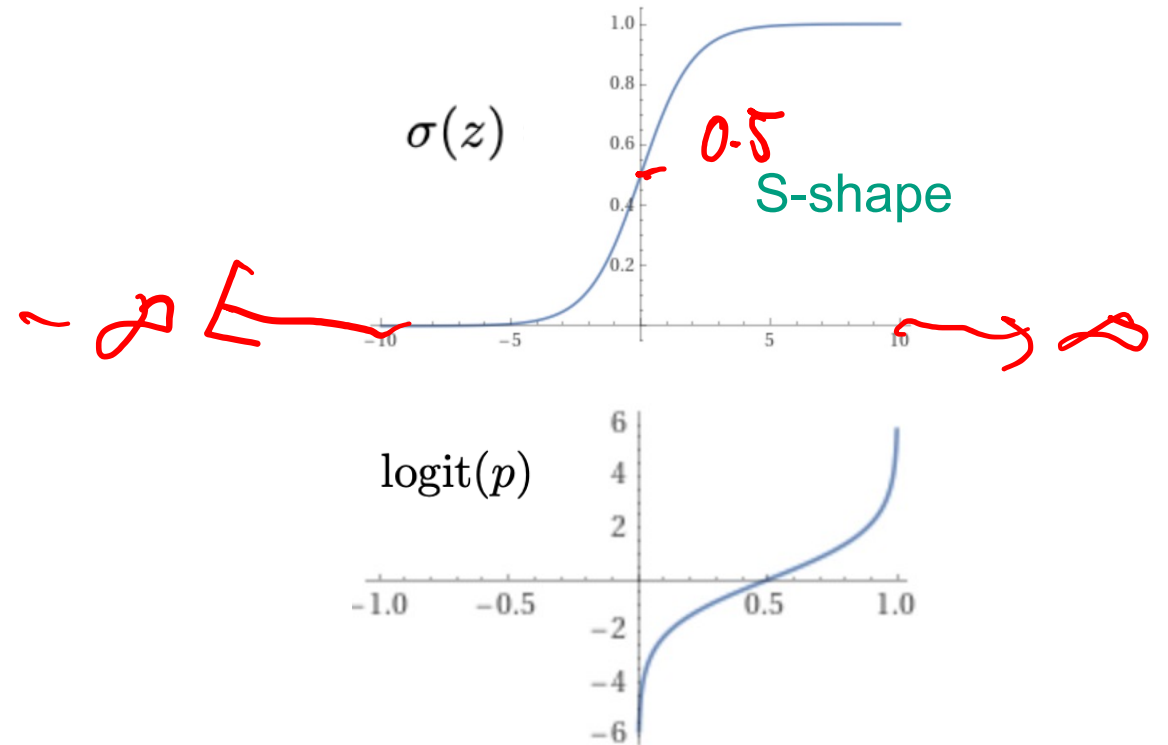
But we are interested in the probability

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

Logistic function / sigmoid function

$$\sigma : \mathbb{R} \rightarrow (0, 1), \quad \sigma(z) = \frac{1}{1 + \exp(-z)}$$

Sigmoid is the inverse of logit





Logistic regression, the plan

Use linear model for the log-odds

Convert to probability using logistic function

Use thresholding to predict class label

$$\hat{y}(z) = \begin{cases} 1 & \text{if } \sigma(z) \leq 0.5 \\ 0 & \text{otherwise} \end{cases} = \begin{cases} 1 & \text{if } z \leq 0 \\ 0 & \text{otherwise} \end{cases}.$$

Logistic regression

$$z = \mathbf{x}^{(i)} \cdot \boldsymbol{\theta}$$

Probabilities of the data

$$P(y^{(i)} = 1 | \mathbf{x}^{(i)}; \boldsymbol{\theta}) = \sigma(z) \quad P(y^{(i)} = 0 | \mathbf{x}^{(i)}; \boldsymbol{\theta}) = 1 - \sigma(z)$$

Combined

$= p$

$= 1 - p$

$$P(y^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta}) = \sigma(z)^{y^{(i)}} (1 - \sigma(z))^{(1-y^{(i)})}$$

$$y^{(i)} \in \{0, 1\}$$

Like for linear regression: maximize likelihood of the parameters

$$\mathcal{L}(\boldsymbol{\theta}) = P(\mathbf{y} | \mathbf{X}; \boldsymbol{\theta})$$

$\sigma(z)$
 $1 - \sigma(z)$

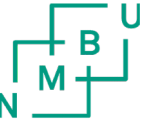




Logistic regression

Goal: maximize likelihood of the parameters

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}) &= P(\mathbf{y} \mid \mathbf{X}; \boldsymbol{\theta}) \\ &= \prod_{i=1}^n P(y^{(i)} \mid x^{(i)}; \boldsymbol{\theta}), \quad (\text{samples are i.i.d.}) \\ &= \prod_{i=1}^n \left[\sigma(z)^{y^{(i)}} (1 - \sigma(z))^{(1-y^{(i)})} \right].\end{aligned}$$

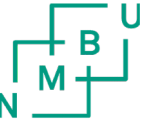


Logistic regression

Goal: maximize likelihood of the parameters

Use log-likelihood instead

$$\begin{aligned}\ell(\boldsymbol{\theta}) &= \log \mathcal{L}(\boldsymbol{\theta}) = \log P(\mathbf{y} \mid \mathbf{X}; \boldsymbol{\theta}) \\ &= \log \prod_{i=1}^n \left[\sigma(z)^{y^{(i)}} (1 - \sigma(z))^{(1-y^{(i)})} \right] \\ &= \sum_{i=1}^n \left[y^{(i)} \log(\sigma(z)) + (1 - y^{(i)}) \log(1 - \sigma(z)) \right]\end{aligned}$$



Logistic regression

Goal: maximize likelihood of the parameters

Use log-likelihood instead

$$\begin{aligned}\ell(\boldsymbol{\theta}) &= \log \mathcal{L}(\boldsymbol{\theta}) = \log P(\mathbf{y} \mid \mathbf{X}; \boldsymbol{\theta}) \\ &= \log \prod_{i=1}^n \left[\sigma(z)^{y^{(i)}} (1 - \sigma(z))^{(1-y^{(i)})} \right] \\ &= \sum_{i=1}^n \left[y^{(i)} \log(\sigma(z)) + (1 - y^{(i)}) \log(1 - \sigma(z)) \right]\end{aligned}$$

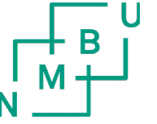
Logistic regression

Goal: maximize likelihood of the parameters → minimize negative log-likelihood

Use log-likelihood instead

$$L(\boldsymbol{\theta}) = -\ell(\boldsymbol{\theta}) \quad \ell(\boldsymbol{\theta}) = \sum_{i=1}^n \left[y^{(i)} \log(\sigma(z)) + (1 - y^{(i)}) \log(1 - \sigma(z)) \right]$$

$$\frac{\partial L(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathbf{0} \quad \text{condition for minimum}$$



Logistic regression

$$\sum_{i=1}^n \left[y^{(i)} \log(\sigma(z)) + (1 - y^{(i)}) \log(1 - \sigma(z)) \right]$$

$$\frac{\partial \sigma(z)}{\partial z} =$$

$$\frac{\partial (\log \sigma(z))}{\partial \theta} =$$

$$\frac{\partial (\log(1 - \sigma(z)))}{\partial \theta} =$$

See lecture notes
PDF

Logistic regression

Goal: maximize likelihood of the parameters → minimize negative log-likelihood

Use log-likelihood instead

$$L(\boldsymbol{\theta}) = -\ell(\boldsymbol{\theta}) \quad \ell(\boldsymbol{\theta}) = \sum_{i=1}^n \left[y^{(i)} \log(\sigma(z)) + (1 - y^{(i)}) \log(1 - \sigma(z)) \right]$$

$$\begin{aligned} \frac{\partial L(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} &= - \sum_{i=1}^n \left[y^{(i)} (1 - \sigma(z)) \mathbf{x}^{(i)} - (1 - y^{(i)}) \sigma(z) \mathbf{x}^{(i)} \right] \\ &= - \sum_{i=1}^n \left[(y^{(i)} - \sigma(z)) \mathbf{x}^{(i)} \right] = \mathbf{0}, \end{aligned}$$

Component-wise

$$\frac{\partial L(\boldsymbol{\theta})}{\partial \theta_j} = - \sum_{i=1}^n \left(x_j^{(i)} \left[(y^{(i)} - \sigma(\mathbf{x}^{(i)} \boldsymbol{\theta})) \right] \right) = 0$$

Handwritten red annotations: An arrow points from the term $(y^{(i)} - \sigma(\mathbf{x}^{(i)} \boldsymbol{\theta}))$ in the equation to the handwritten expression $x^{(i)} \theta$. A red underline is drawn under the entire equation.

Same gradient as for linear regression / Adaline, except for $\sigma(z)$!

But nonlinear, no explicit solution!

Logistic regression summary

Learn / train / fit

$$\frac{\partial L(\boldsymbol{\theta})}{\partial \theta_j} = - \sum_{i=1}^n \left(x_j^{(i)} \left[(y^{(i)} - \sigma(\mathbf{x}^{(i)} \boldsymbol{\theta})) \right] \right) = 0.$$

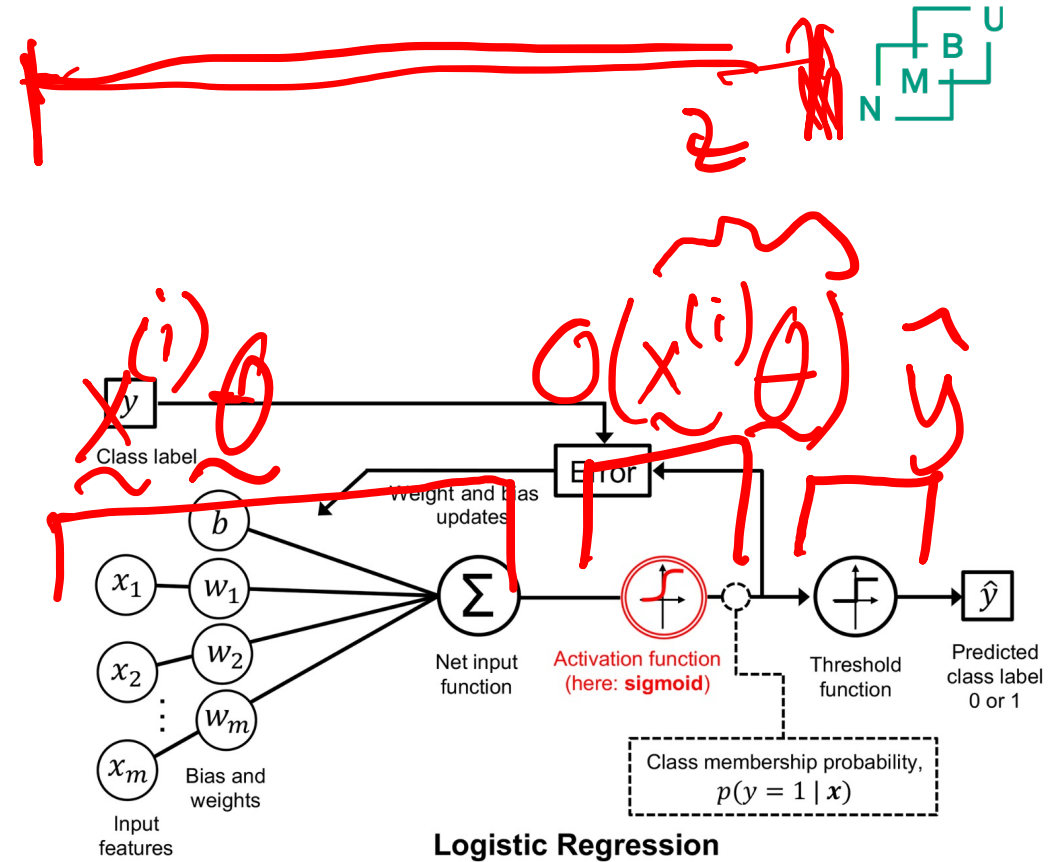
$$\boldsymbol{\theta}^{n+1} = \boldsymbol{\theta} - \eta \frac{\partial L(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \quad \text{(Batch) gradient descent}$$

Predict

$$z = \mathbf{X}\boldsymbol{\theta} \quad \mathbf{y} = \sigma(z) = \frac{1}{1 + \exp(-z)}$$

→ log odds

$\hat{y} = \text{threshold}(y)$

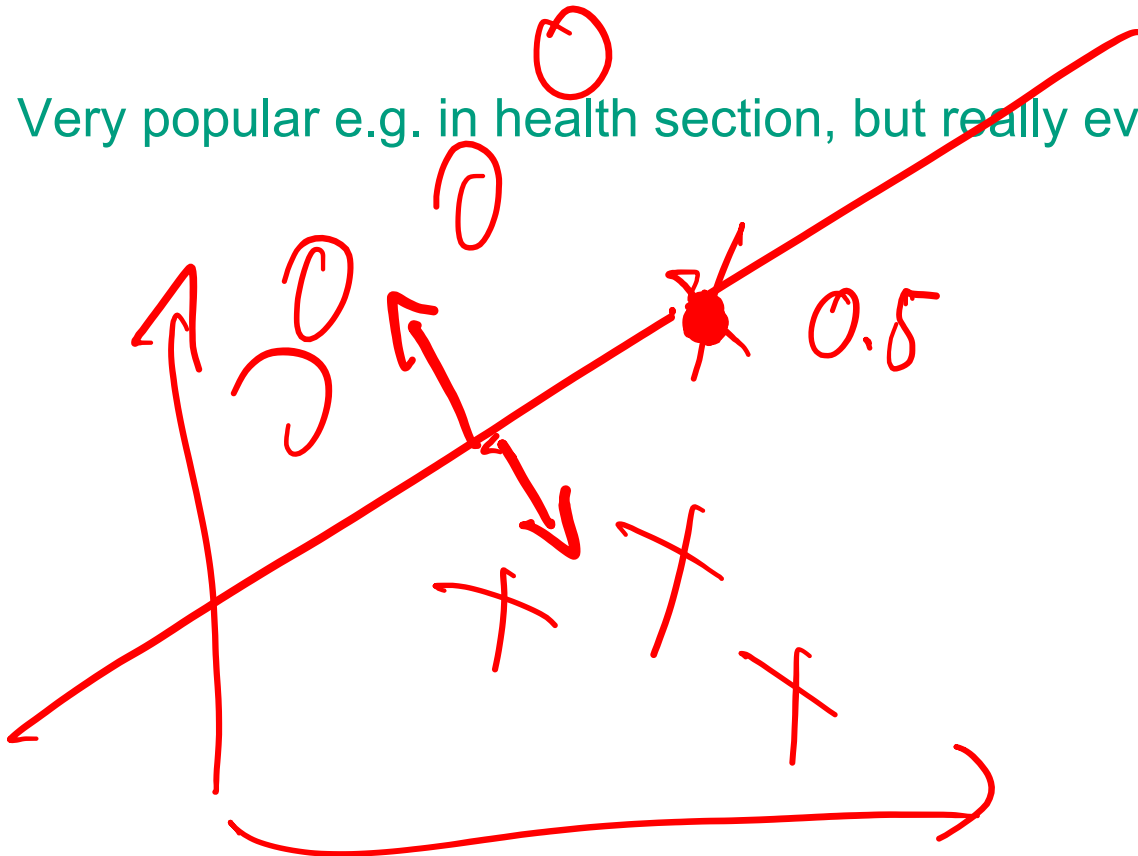


log_regression.ipynb

Logistic regression summary

Logistic regression gives us label **and** probability

Very popular e.g. in health section, but really everywhere



`log_regression.ipynb`

