

1 Linear regression

- $m \in \mathbb{N}$, number of examples.
- $n \in \mathbb{N}$, example size.
- $A \bowtie B$, element-wise using broadcasting (python like).

$$\mathbf{X} \in \mathbb{R}^{m \times n}$$

$$\mathbf{x}_i \in \mathbb{R}^{1 \times n}$$

$$\mathbf{y} \in \mathbb{R}^{m \times 1}$$

$$\mathbf{w} \in \mathbb{R}^{1 \times n}$$

$$b \in \mathbb{R}$$

1.1 Cost function

$$J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) = \frac{1}{m} \sum_{i=1}^m ([\mathbf{w}\mathbf{x}_i^\top + b] - y_i)^2$$

$$J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) = \frac{1}{m} \sum_{i=1}^m (\mathbf{w}\mathbf{x}_i^\top \mathbf{x}_i \mathbf{w}^\top + 2b\mathbf{w}\mathbf{x}_i^\top + b^2 - 2\mathbf{w}\mathbf{x}_i^\top y_i - 2by_i + y_i^2)$$

$$J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) = \mathbf{w} \left[\frac{1}{m} \sum_{i=1}^m \mathbf{x}_i^\top \mathbf{x}_i \right] \mathbf{w}^\top + 2b\mathbf{w} \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i^\top + \frac{1}{m} \sum_{i=1}^m b^2 - 2\mathbf{w} \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i^\top y_i - 2b \frac{1}{m} \sum_{i=1}^m y_i + \frac{1}{m} \sum_{i=1}^m y_i^2$$

$$J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) = \mathbf{w} \left[\frac{1}{m} \mathbf{X}^\top \mathbf{X} \right] \mathbf{w}^\top + 2b\mathbf{w}\boldsymbol{\mu}_2^\top(\mathbf{X}) + b^2 - 2\mathbf{w}\boldsymbol{\mu}_2^\top(\mathbf{X} \bowtie \mathbf{y}) - 2b\boldsymbol{\mu}(\mathbf{y}) + \frac{1}{m} \mathbf{y}^\top \mathbf{y}$$

1.2 Gradient (weights)

$$\frac{\partial}{\partial w^{(j)}} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w^{(j)}} ([\mathbf{w}\mathbf{x}_i^\top + b] - y_i)^2$$

$$\frac{\partial}{\partial w^{(j)}} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) = \frac{2}{m} \sum_{i=1}^m ([\mathbf{w}\mathbf{x}_i^\top + b] - y_i) \frac{\partial}{\partial w^{(j)}} ([\mathbf{w}\mathbf{x}_i^\top + b] - y_i)$$

$$\frac{\partial}{\partial w^{(j)}} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) = \frac{2}{m} \sum_{i=1}^m ([\mathbf{w}\mathbf{x}_i^\top + b] - y_i) \left(\frac{\partial}{\partial w^{(j)}} \mathbf{w}\mathbf{x}_i^\top + \frac{\partial}{\partial w^{(j)}} b - \frac{\partial}{\partial w^{(j)}} y_i \right)$$

$$\frac{\partial}{\partial w^{(j)}} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) = \frac{2}{m} \sum_{i=1}^m ([\mathbf{w}\mathbf{x}_i^\top + b] - y_i) \left(\frac{\partial}{\partial w^{(j)}} \mathbf{w}\mathbf{x}_i^\top \right)$$

$$\begin{aligned}
\frac{\partial}{\partial w^{(j)}} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= \frac{2}{m} \sum_{i=1}^m (\mathbf{w} \mathbf{x}_i^\top + b - y_i) x_i^{(j)} \\
\frac{\partial}{\partial w^{(j)}} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= \mathbf{w} \frac{2}{m} \sum_{i=1}^m \mathbf{x}_i^\top x_i^{(j)} + b \frac{2}{m} \sum_{i=1}^m x_i^{(j)} - \frac{2}{m} \sum_{i=1}^m y_i x_i^{(j)} \\
\frac{\partial}{\partial \mathbf{w}} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= \frac{2}{m} \begin{bmatrix} \mathbf{w} \sum_{i=1}^m \mathbf{x}_i^\top x_i^{(1)} \\ \mathbf{w} \sum_{i=1}^m \mathbf{x}_i^\top x_i^{(2)} \\ \vdots \\ \mathbf{w} \sum_{i=1}^m \mathbf{x}_i^\top x_i^{(n)} \end{bmatrix} + \frac{2b}{m} \begin{bmatrix} \sum_{i=1}^m x_i^{(1)} \\ \sum_{i=1}^m x_i^{(2)} \\ \vdots \\ \sum_{i=1}^m x_i^{(n)} \end{bmatrix} - \frac{2}{m} \begin{bmatrix} \sum_{i=1}^m y_i x_i^{(1)} \\ \sum_{i=1}^m y_i x_i^{(2)} \\ \vdots \\ \sum_{i=1}^m y_i x_i^{(n)} \end{bmatrix} \\
\frac{\partial}{\partial \mathbf{w}} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= \frac{2}{m} \mathbf{w} \sum_{i=1}^m \mathbf{x}_i^\top \mathbf{x}_i + \frac{2b}{m} \sum_{i=1}^m \mathbf{x}_i - \frac{2}{m} \sum_{i=1}^m y_i \mathbf{x}_i \\
\boxed{\frac{\partial}{\partial \mathbf{w}} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= 2 \left(\mathbf{w} \left[\frac{1}{m} \mathbf{X}^\top \mathbf{X} \right] + b \boldsymbol{\mu}_2(\mathbf{X}) - \boldsymbol{\mu}_2(\mathbf{X} \bowtie \mathbf{y}) \right)}
\end{aligned}$$

1.3 Gradient (bias)

$$\begin{aligned}
\frac{\partial}{\partial b} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial b} ([\mathbf{w} \mathbf{x}_i^\top + b] - y_i)^2 \\
\frac{\partial}{\partial b} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= \frac{2}{m} \sum_{i=1}^m ([\mathbf{w} \mathbf{x}_i^\top + b] - y_i) \frac{\partial}{\partial b} ([\mathbf{w} \mathbf{x}_i^\top + b] - y_i) \\
\frac{\partial}{\partial b} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= \frac{2}{m} \sum_{i=1}^m ([\mathbf{w} \mathbf{x}_i^\top + b] - y_i) \left(\frac{\partial}{\partial b} \mathbf{w} \mathbf{x}_i^\top + \frac{\partial}{\partial b} b - \frac{\partial}{\partial b} y_i \right) \\
\frac{\partial}{\partial b} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= \frac{2}{m} \sum_{i=1}^m ([\mathbf{w} \mathbf{x}_i^\top + b] - y_i) \left(\frac{\partial}{\partial b} b \right) \\
\frac{\partial}{\partial b} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= \frac{2}{m} \sum_{i=1}^m (\mathbf{w} \mathbf{x}_i^\top + b - y_i) \\
\frac{\partial}{\partial b} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= \mathbf{w} \frac{2}{m} \sum_{i=1}^m \mathbf{x}_i^\top + \frac{2}{m} \sum_{i=1}^m b - \frac{2}{m} \sum_{i=1}^m y_i \\
\boxed{\frac{\partial}{\partial b} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= 2 (\mathbf{w} \boldsymbol{\mu}_2^\top(\mathbf{X}) + b - \mu(\mathbf{y}))}
\end{aligned}$$

1.4 Analitic solution

$$\frac{\partial}{\partial \mathbf{w}} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) = \mathbf{0} = 2 \left(\mathbf{w} \left[\frac{1}{m} \mathbf{X}^\top \mathbf{X} \right] + b \boldsymbol{\mu}_2(\mathbf{X}) - \boldsymbol{\mu}_2(\mathbf{X} \bowtie \mathbf{y}) \right)$$

$$\frac{\partial}{\partial b} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) = 0 = 2 (\mathbf{w} \boldsymbol{\mu}_2^\top(\mathbf{X}) + b - \mu(\mathbf{y}))$$

$$\begin{bmatrix} \mathbf{w} \left[\frac{1}{m} \mathbf{X}^\top \mathbf{X} \right] + b \boldsymbol{\mu}_2(\mathbf{X}) - \boldsymbol{\mu}_2(\mathbf{X} \bowtie \mathbf{y}) \\ \mathbf{w} \boldsymbol{\mu}_2^\top(\mathbf{X}) + b - \mu(\mathbf{y}) \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} \mathbf{w} \left[\frac{1}{m} \mathbf{X}^\top \mathbf{X} \right] + b \boldsymbol{\mu}_2(\mathbf{X}) \\ \mathbf{w} \boldsymbol{\mu}_1^\top(\mathbf{X}) + b \end{bmatrix} = \begin{bmatrix} \boldsymbol{\mu}_2(\mathbf{X} \bowtie \mathbf{y}) \\ \mu(\mathbf{y}) \end{bmatrix}$$

$$\begin{bmatrix} \mathbf{w} & b \end{bmatrix} \begin{bmatrix} \left[\frac{1}{m} \mathbf{X}^\top \mathbf{X} \right] & \boldsymbol{\mu}_2^\top(\mathbf{X}) \\ \boldsymbol{\mu}_2(\mathbf{X}) & 1 \end{bmatrix} = \begin{bmatrix} \boldsymbol{\mu}_2(\mathbf{X} \bowtie \mathbf{y}) \\ \mu(\mathbf{y}) \end{bmatrix}$$

$$\boxed{\begin{bmatrix} \mathbf{w} & b \end{bmatrix} = \begin{bmatrix} \boldsymbol{\mu}_2(\mathbf{X} \bowtie \mathbf{y}) \\ \mu(\mathbf{y}) \end{bmatrix} \begin{bmatrix} \left[\frac{1}{m} \mathbf{X}^\top \mathbf{X} \right] & \boldsymbol{\mu}_2^\top(\mathbf{X}) \\ \boldsymbol{\mu}_2(\mathbf{X}) & 1 \end{bmatrix}^{-1}}$$

2 Logistic regression

- $m \in \mathbb{N}$, number of examples.
- $n \in \mathbb{N}$, example size.

$$\mathbf{X} \in \mathbb{R}^{m \times n}$$

$$\mathbf{x}_i \in \mathbb{R}^{1 \times n}$$

$$\mathbf{y} \in \mathbb{B}^{m \times 1}$$

$$\mathbf{w} \in \mathbb{R}^{1 \times n}$$

$$b \in \mathbb{R}$$

2.1 Cost function

$$J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) = -\frac{1}{m} \sum_{i=1}^m (y_i \log(\sigma(\mathbf{w} \mathbf{x}_i^\top + b)) + (1 - y_i) \log(1 - \sigma(\mathbf{w} \mathbf{x}_i^\top + b)))$$

$$\boxed{J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) = -\frac{1}{m} \left[\sum_{i=1}^{m \wedge y_i=1} \log(\sigma(\mathbf{w} \mathbf{x}_i^\top + b)) + \sum_{i=1}^{m \wedge y_i=0} \log(1 - \sigma(\mathbf{w} \mathbf{x}_i^\top + b)) \right]}$$

2.2 Gradient (weights)

$$\begin{aligned}
\frac{\partial}{\partial w^{(j)}} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= -\frac{1}{m} \left[\sum_{i=1}^{m \wedge y_i=1} \frac{\partial}{\partial w^{(j)}} \log(\sigma(\mathbf{w}\mathbf{x}_i^\top + b)) + \sum_{i=1}^{m \wedge y_i=0} \frac{\partial}{\partial w^{(j)}} \log(1 - \sigma(\mathbf{w}\mathbf{x}_i^\top + b)) \right] \\
\frac{\partial}{\partial w^{(j)}} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= -\frac{1}{m} \left[\sum_{i=1}^{m \wedge y_i=1} \frac{1}{\sigma(\mathbf{w}\mathbf{x}_i^\top + b)} \frac{\partial}{\partial w^{(j)}} \sigma(\mathbf{w}\mathbf{x}_i^\top + b) + \sum_{i=1}^{m \wedge y_i=0} \frac{\partial}{\partial w^{(j)}} \log(1 - \sigma(\mathbf{w}\mathbf{x}_i^\top + b)) \right] \\
\frac{\partial}{\partial w^{(j)}} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= -\frac{1}{m} \left[\sum_{i=1}^{m \wedge y_i=1} \frac{1}{\cancel{\sigma(\mathbf{w}\mathbf{x}_i^\top + b)}} \left(\cancel{\sigma(\mathbf{w}\mathbf{x}_i^\top + b)} (1 - \sigma(\mathbf{w}\mathbf{x}_i^\top + b)) \right) \frac{\partial}{\partial w^{(j)}} (\mathbf{w}\mathbf{x}_i^\top + b) + \sum_{i=1}^{m \wedge y_i=0} \frac{\partial}{\partial w^{(j)}} \log(1 - \sigma(\mathbf{w}\mathbf{x}_i^\top + b)) \right] \\
\frac{\partial}{\partial w^{(j)}} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= -\frac{1}{m} \left[\sum_{i=1}^{m \wedge y_i=1} \left(x_i^{(j)} - \sigma(\mathbf{w}\mathbf{x}_i^\top + b) x_i^{(j)} \right) + \sum_{i=1}^{m \wedge y_i=0} \frac{\partial}{\partial w^{(j)}} \log(1 - \sigma(\mathbf{w}\mathbf{x}_i^\top + b)) \right] \\
\frac{\partial}{\partial w^{(j)}} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= -\frac{1}{m} \left[\sum_{i=1}^{m \wedge y_i=1} \left(x_i^{(j)} - \sigma(\mathbf{w}\mathbf{x}_i^\top + b) x_i^{(j)} \right) + \sum_{i=1}^{m \wedge y_i=0} \frac{1}{1 - \sigma(\mathbf{w}\mathbf{x}_i^\top + b)} \frac{\partial}{\partial w^{(j)}} (1 - \sigma(\mathbf{w}\mathbf{x}_i^\top + b)) \right] \\
\frac{\partial}{\partial w^{(j)}} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= -\frac{1}{m} \left[\sum_{i=1}^{m \wedge y_i=1} \left(x_i^{(j)} - \sigma(\mathbf{w}\mathbf{x}_i^\top + b) x_i^{(j)} \right) - \sum_{i=1}^{m \wedge y_i=0} \frac{1}{1 - \sigma(\mathbf{w}\mathbf{x}_i^\top + b)} \left(\frac{\partial}{\partial w^{(j)}} \sigma(\mathbf{w}\mathbf{x}_i^\top + b) \right) \right] \\
\frac{\partial}{\partial w^{(j)}} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= -\frac{1}{m} \left[\sum_{i=1}^{m \wedge y_i=1} \left(x_i^{(j)} - \sigma(\mathbf{w}\mathbf{x}_i^\top + b) x_i^{(j)} \right) - \sum_{i=1}^{m \wedge y_i=0} \frac{1}{\cancel{1 - \sigma(\mathbf{w}\mathbf{x}_i^\top + b)}} \left(\cancel{\sigma(\mathbf{w}\mathbf{x}_i^\top + b)} (1 - \cancel{\sigma(\mathbf{w}\mathbf{x}_i^\top + b)}) \frac{\partial}{\partial w^{(j)}} (\mathbf{w}\mathbf{x}_i^\top + b) \right) \right] \\
\frac{\partial}{\partial w^{(j)}} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= -\frac{1}{m} \left[\sum_{i=1}^{m \wedge y_i=1} \left(x_i^{(j)} - \sigma(\mathbf{w}\mathbf{x}_i^\top + b) x_i^{(j)} \right) - \sum_{i=1}^{m \wedge y_i=0} \sigma(\mathbf{w}\mathbf{x}_i^\top + b) x_i^{(j)} \right] \\
\frac{\partial}{\partial w^{(j)}} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= -\frac{1}{m} \left[\sum_{i=1}^m y_i \left(x_i^{(j)} - \sigma(\mathbf{w}\mathbf{x}_i^\top + b) x_i^{(j)} \right) - \sum_{i=1}^m (1 - y_i) \sigma(\mathbf{w}\mathbf{x}_i^\top + b) x_i^{(j)} \right] \\
\frac{\partial}{\partial w^{(j)}} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= -\frac{1}{m} \sum_{i=1}^m \left[x_i^{(j)} y_i \cancel{- \sigma(\mathbf{w}\mathbf{x}_i^\top + b) x_i^{(j)}} + \cancel{\sigma(\mathbf{w}\mathbf{x}_i^\top + b) x_i^{(j)}} y_i - \sigma(\mathbf{w}\mathbf{x}_i^\top + b) x_i^{(j)} \right] \\
\frac{\partial}{\partial w^{(j)}} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= \frac{1}{m} \sum_{i=1}^m \left[\sigma(\mathbf{w}\mathbf{x}_i^\top + b) x_i^{(j)} - x_i^{(j)} y_i \right] \\
\frac{\partial}{\partial w^{(j)}} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= \frac{1}{m} \sum_{i=1}^m x_i^{(j)} \sigma(\mathbf{w}\mathbf{x}_i^\top + b) - \frac{1}{m} \sum_{i=1}^m x_i^{(j)} y_i \\
\frac{\partial}{\partial \mathbf{w}} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= \frac{1}{m} \sum_{i=1}^m x_i^{(j)} \sigma(\mathbf{w}\mathbf{x}_i^\top + b) - \boldsymbol{\mu}_1(\mathbf{X} \bowtie \mathbf{y}) \\
\boxed{\frac{\partial}{\partial \mathbf{w}} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= \boldsymbol{\mu}_1(\mathbf{X} \bowtie \sigma(\mathbf{X}\mathbf{w}^\top + \mathbf{b})) - \boldsymbol{\mu}_1(\mathbf{X} \bowtie \mathbf{y})}
\end{aligned}$$

2.3 Gradient (bias)

$$\begin{aligned}
\frac{\partial}{\partial b} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= -\frac{1}{m} \left[\sum_{i=1}^{m \wedge y_i=1} \frac{\partial}{\partial b} \log(\sigma(\mathbf{w}\mathbf{x}_i^\top + b)) + \sum_{i=1}^{m \wedge y_i=0} \frac{\partial}{\partial b} \log(1 - \sigma(\mathbf{w}\mathbf{x}_i^\top + b)) \right] \\
\frac{\partial}{\partial b} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= -\frac{1}{m} \left[\sum_{i=1}^{m \wedge y_i=1} \frac{1}{\sigma(\mathbf{w}\mathbf{x}_i^\top + b)} \frac{\partial}{\partial b} \sigma(\mathbf{w}\mathbf{x}_i^\top + b) + \sum_{i=1}^{m \wedge y_i=0} \frac{1}{1 - \sigma(\mathbf{w}\mathbf{x}_i^\top + b)} \frac{\partial}{\partial b} (1 - \sigma(\mathbf{w}\mathbf{x}_i^\top + b)) \right] \\
\frac{\partial}{\partial b} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= -\frac{1}{m} \left[\sum_{i=1}^{m \wedge y_i=1} \frac{1}{\cancel{\sigma(\mathbf{w}\mathbf{x}_i^\top + b)}} \left(\cancel{\sigma(\mathbf{w}\mathbf{x}_i^\top + b)} (1 - \sigma(\mathbf{w}\mathbf{x}_i^\top + b)) \right) - \sum_{i=1}^{m \wedge y_i=0} \frac{1}{\cancel{1 - \sigma(\mathbf{w}\mathbf{x}_i^\top + b)}} \left(\sigma(\mathbf{w}\mathbf{x}_i^\top + b) \cancel{(1 - \sigma(\mathbf{w}\mathbf{x}_i^\top + b))} \right) \right] \\
\frac{\partial}{\partial b} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= \frac{1}{m} \left[\sum_{i=1}^{m \wedge y_i=0} \sigma(\mathbf{w}\mathbf{x}_i^\top + b) - \sum_{i=1}^{m \wedge y_i=1} (1 - \sigma(\mathbf{w}\mathbf{x}_i^\top + b)) \right] \\
\frac{\partial}{\partial b} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= \frac{1}{m} \left[\sum_{i=1}^m (1 - y_i) \sigma(\mathbf{w}\mathbf{x}_i^\top + b) - y_i (1 - \sigma(\mathbf{w}\mathbf{x}_i^\top + b)) \right] \\
\frac{\partial}{\partial b} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= \frac{1}{m} \left[\sum_{i=1}^m \sigma(\mathbf{w}\mathbf{x}_i^\top + b) \cancel{-y_i \sigma(\mathbf{w}\mathbf{x}_i^\top + b)} - y_i \cancel{+y_i \sigma(\mathbf{w}\mathbf{x}_i^\top + b)} \right] \\
\boxed{\frac{\partial}{\partial b} J(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) &= \frac{1}{m} \sum_{i=1}^m \sigma(\mathbf{w}\mathbf{x}_i^\top + b) - \mu(\mathbf{y})}
\end{aligned}$$