

Neural Computation

28 September 2023

Linear regression: toy example

- Commute time on bus
 - Want to predict commute time to University
 - Input variables (features)?
 - Distance to University
 - Day of the week
 - Output / target?
 - Commute time
 - Data

$$f(x) = x w$$

| Dist (km) | Commute time (min) |
|-----------|--------------------|
| 2.7 | 25 |
| 4.1 | 33 |
| 1.0 | 15 |
| 5.2 | 45 |
| 2.8 | 22 |



Linear regression: least squares (1d)

- Let $d = 1$, then:

$$C(w) = \frac{1}{2n} \sum_{i=1}^n (y^i - x^i w)^2$$

How important is this factor?

| i | x | y |
|-----|--------------|-----------------------|
| | Dist (km) | Commute time (min) |
| | 2.7 | 25 |
| | 4.1 | 33 |
| | 1.0 | 15 |
| | 5.2 | 45 |
| | 2.8 | 22 |

Linear regression: least squares (1d)

- Let $d = 1$, then:

$$C(w) = \frac{1}{2n} \sum_{i=1}^n (y^i - x^i w)^2$$

What shape does this function have?

| i | x | y |
|-----|-----------|--------------------|
| | Dist (km) | Commute time (min) |
| | 2.7 | 25 |
| | 4.1 | 33 |
| | 1.0 | 15 |
| | 5.2 | 45 |
| | 2.8 | 22 |

Linear regression: least squares (1d)

- Let $d = 1$, then:

$$C(w) = \frac{1}{2n} \sum_{i=1}^n (y^i - x^i w)^2 = \frac{1}{2n} \sum_{i=1}^n \left(\underbrace{x^{i^2} w^2}_{\text{quadratic}} - \underbrace{2y^i x^i w}_{\text{linear}} + \underbrace{y^{i^2}}_{\text{constant}} \right)$$

How can we find the best w ?

| i | x | y |
|-----|--------------|-----------------------|
| | Dist (km) | Commute time (min) |
| | 2.7 | 25 |
| | 4.1 | 33 |
| | 1.0 | 15 |
| | 5.2 | 45 |
| | 2.8 | 22 |

Linear regression: least squares (1d)

- Let $d = 1$, then:

$$C(w) = \frac{1}{2n} \sum_{i=1}^n (y^i - x^i w)^2 = \frac{1}{2n} \sum_{i=1}^n \left(\underbrace{x^{i^2} w^2}_{\text{quadratic}} - \underbrace{2y^i x^i w}_{\text{linear}} + \underbrace{y^{i^2}}_{\text{constant}} \right)$$

- It then follows that:

$$C'(w) = \frac{1}{2n} \sum_{i=1}^n (2x^{i^2} w - 2y^i x^i) = \frac{1}{n} \sum_{i=1}^n x^{i^2} w - \frac{1}{n} \sum_{i=1}^n y^i x^i$$

| i | x | y |
|-----|-----------|--------------------|
| | Dist (km) | Commute time (min) |
| | 2.7 | 25 |
| | 4.1 | 33 |
| | 1.0 | 15 |
| | 5.2 | 45 |
| | 2.8 | 22 |

Linear regression: least squares (1d)

- Let $d = 1$, then:

$$C(w) = \frac{1}{2n} \sum_{i=1}^n (y^i - x^i w)^2 = \frac{1}{2n} \sum_{i=1}^n \left(\underbrace{x^{i2} w^2}_{\text{quadratic}} - \underbrace{2y^i x^i w}_{\text{linear}} + \underbrace{y^{i2}}_{\text{constant}} \right)$$

- It then follows that:

$$C'(w) = \frac{1}{2n} \sum_{i=1}^n (2x^{i2} w - 2y^i x^i) = \frac{1}{n} \sum_{i=1}^n x^{i2} w - \frac{1}{n} \sum_{i=1}^n y^i x^i$$

- According to the first-order optimality condition, we know the optimal w^* satisfies

$$C'(w^*) = 0 \Rightarrow \frac{1}{n} \sum_{i=1}^n x^{i2} w^* = \frac{1}{n} \sum_{i=1}^n y^i x^i$$

- It then follows that:

$$w^* = \frac{\sum_{i=1}^n y^i x^i}{\sum_{i=1}^n x^{i2}}$$

Linear regression: toy example

- Commute time on bus
 - Want to predict commute time to University
 - Input variables (features)?
 - Distance to University
 - Day of the week
 - Output / target?
 - Commute time
 - Data
 - day = 1 if weekday, day = 0 otherwise

| Dist (km) | Day | Commute time (min) |
|-----------|-----|--------------------|
| 2.7 | 1 | 25 |
| 4.1 | 1 | 33 |
| 1.0 | 0 | 15 |
| 5.2 | 1 | 45 |
| 2.8 | 0 | 22 |



$$f(\mathbf{x}) = w_1x_1 + \dots + w_dx_d = (w_0 + w_1 + \dots + w_d) \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix} = \mathbf{w}^T \mathbf{x}$$

Linear regression: matrix form

- Let's recall that $C(\mathbf{w}) = \frac{1}{2n} \sum_{i=1}^n (y^i - \mathbf{x}^{iT} \mathbf{w})^2$
- Let's consider $\mathbf{x}^{iT} = (x_1^i, x_2^i, \dots, x_d^i)$

$$X = \begin{pmatrix} \mathbf{x}^{1T} \\ \vdots \\ \mathbf{x}^{nT} \end{pmatrix} \in \mathbb{R}^{n \times d}, \mathbf{y} = \begin{pmatrix} y^1 \\ \vdots \\ y^n \end{pmatrix} \in \mathbb{R}^n \Rightarrow X\mathbf{w} - \mathbf{y} = \begin{pmatrix} \mathbf{x}^{1T} \mathbf{w} - y^1 \\ \vdots \\ \mathbf{x}^{nT} \mathbf{w} - y^n \end{pmatrix}$$

| Dist (km) | Day | Commute time (min) |
|-----------|-------|--------------------|
| x_1 | x_2 | y |
| 2.7 | 1 | 25 |
| 4.1 | 1 | 33 |
| 1.0 | 0 | 15 |
| 5.2 | 1 | 45 |
| 2.8 | 0 | 22 |

$$\mathbf{y} = \begin{pmatrix} 25 \\ 33 \\ 15 \\ 45 \\ 22 \end{pmatrix}, X = \begin{pmatrix} 2.7 & 1 \\ 4.1 & 1 \\ 1.0 & 0 \\ 5.2 & 1 \\ 2.8 & 0 \end{pmatrix}, \mathbf{w} = \begin{pmatrix} w_0 \\ w_1 \\ w_2 \end{pmatrix}$$

Linear regression: matrix form

- Let's recall that $C(\mathbf{w}) = \frac{1}{2n} \sum_{i=1}^n (y^i - \mathbf{x}^{iT} \mathbf{w})^2$
- Let's consider $\mathbf{x}^{iT} = (x_1^i, x_2^i, \dots, x_d^i)$

$$X = \begin{pmatrix} \mathbf{x}^{1T} \\ \vdots \\ \mathbf{x}^{nT} \end{pmatrix} \in \mathbb{R}^{n \times d}, \mathbf{y} = \begin{pmatrix} y^1 \\ \vdots \\ y^n \end{pmatrix} \in \mathbb{R}^n \Rightarrow X\mathbf{w} - \mathbf{y} = \begin{pmatrix} \mathbf{x}^{1T} \mathbf{w} - y^1 \\ \vdots \\ \mathbf{x}^{nT} \mathbf{w} - y^n \end{pmatrix}$$

$$\hat{\mathbf{y}} = X\mathbf{w}$$

$$\begin{aligned} C(\mathbf{w}) &= \frac{1}{2n} (X\mathbf{w} - \mathbf{y})^T (X\mathbf{w} - \mathbf{y}) = \frac{1}{2n} (\mathbf{w}^T X^T - \mathbf{y}^T) (X\mathbf{w} - \mathbf{y}) \\ &= \frac{1}{2n} (\mathbf{w}^T X^T X \mathbf{w} - 2\mathbf{w}^T X^T \mathbf{y} + \mathbf{y}^T \mathbf{y}) \end{aligned}$$

| Dist (km) | Day | Commute time (min) |
|-----------|-------|--------------------|
| x_1 | x_2 | y |
| 2.7 | 1 | 25 |
| 4.1 | 1 | 33 |
| 1.0 | 0 | 15 |
| 5.2 | 1 | 45 |
| 2.8 | 0 | 22 |

$$\mathbf{y} = \begin{pmatrix} 25 \\ 33 \\ 15 \\ 45 \\ 22 \end{pmatrix}, X = \begin{pmatrix} 2.7 & 1 \\ 4.1 & 1 \\ 1.0 & 0 \\ 5.2 & 1 \\ 2.8 & 0 \end{pmatrix}, \mathbf{w} = \begin{pmatrix} w_0 \\ w_1 \\ w_2 \end{pmatrix}$$

Linear regression: closed-form solution

- Let's recall the objective function:

$$C(\mathbf{w}) = \frac{1}{2n} \left(\underbrace{\mathbf{w}^T X^T X \mathbf{w}}_{\text{quadratic}} - \underbrace{2\mathbf{w}^T X^T \mathbf{y}}_{\text{linear}} + \underbrace{\mathbf{y}^T \mathbf{y}}_{\text{constant}} \right)$$

- The gradient of $C(\mathbf{w})$ is:

$$\nabla C(\mathbf{w}) = \frac{1}{2n} (2X^T X \mathbf{w} - 2X^T \mathbf{y})$$

- By the first-order optimality condition, we know that optimal \mathbf{w}^* satisfies:

$$\nabla C(\mathbf{w}^*) = \frac{1}{n} (X^T X \mathbf{w}^* - X^T \mathbf{y}) = 0 \implies X^T X \mathbf{w}^* = X^T \mathbf{y}$$

- If $X^T X$ is invertible, we get $\mathbf{w}^* = (X^T X)^{-1} X^T \mathbf{y}$

Linear regression

- Adding a feature for bias term
 - We can add 1 to get an expanded feature vector

| Dist (km) | Day | Commute time (min) |
|-----------|-------|--------------------|
| x_1 | x_2 | y |
| 2.7 | 1 | 25 |
| 4.1 | 1 | 33 |
| 1.0 | 0 | 15 |
| 5.2 | 1 | 45 |
| 2.8 | 0 | 22 |

\Rightarrow

| One | Dist (km) | Day | Commute time (min) |
|-------|-----------|-------|--------------------|
| x_0 | x_1 | x_2 | y |
| 1 | 2.7 | 1 | 25 |
| 1 | 4.1 | 1 | 33 |
| 1 | 1.0 | 0 | 15 |
| 1 | 5.2 | 1 | 45 |
| 1 | 2.8 | 0 | 22 |

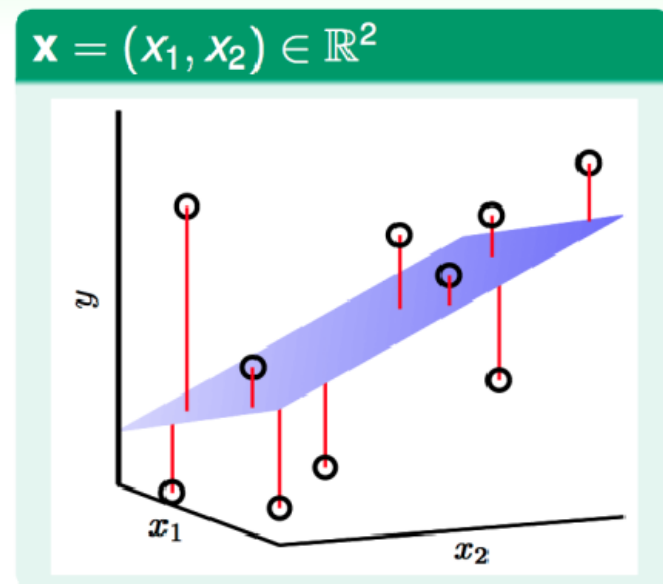
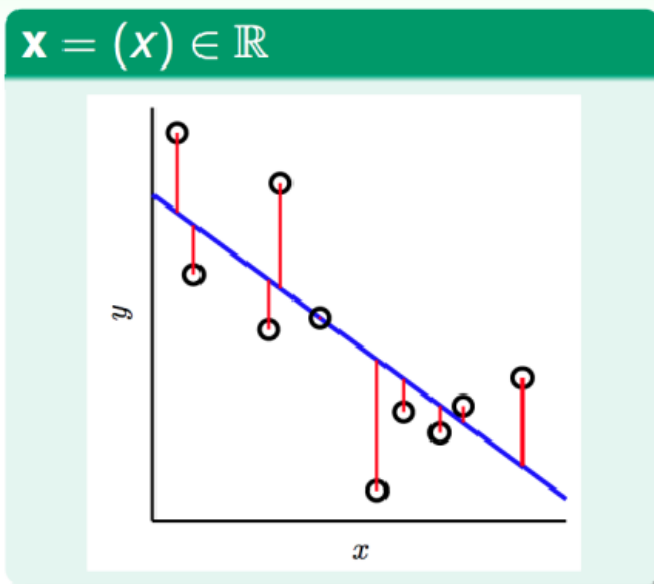
- This allows us to consider the bias in the linear model:

$$f(\mathbf{x}) = w_0 + w_1x_1 + \cdots + w_dx_d = (w_0 + w_1 + \cdots + w_d) \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix} = \mathbf{w}^T \bar{\mathbf{x}}$$

- For brevity, we use notation \mathbf{x} to represent the extended feature $\bar{\mathbf{x}}$:

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$

Linear regression: illustration



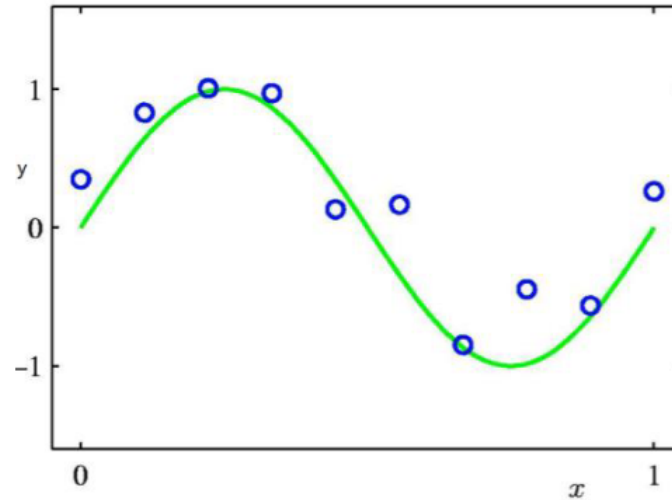
Summary: linear regression

- Linear regression (or least square regression)
 - model linear relationship between input and output (task T)
 - Example points (experience E)
 - mean square error as loss function (performance P)
 - closed-form solution (or exact solution)
 - Add 1s-feature to allow for bias

| One | Dist (km) | Day | Commute time (min) |
|-------|-----------|-------|--------------------|
| x_0 | x_1 | x_2 | y |
| 1 | 2.7 | 1 | 25 |
| 1 | 4.1 | 1 | 33 |
| 1 | 1.0 | 0 | 15 |
| 1 | 5.2 | 1 | 45 |
| 1 | 2.8 | 0 | 22 |

Polynomial regression

- Suppose we want to model the following data



- The input-output relationship is nonlinear!
- How about we try to fit a polynomial?
 - This is known as polynomial regression

$$f(x) = w_0 + w_1x + w_2(x)^2 \dots, w_M(x)^M$$

where $(x)^i$ denotes i^{th} power of x .

- Do we need to derive a whole new regression algorithm?

Remember the 1-feature?

| Dist (km) | Day | Commute time (min) |
|-----------|-------|--------------------|
| x_1 | x_2 | y |
| 2.7 | 1 | 25 |
| 4.1 | 1 | 33 |
| 1.0 | 0 | 15 |
| 5.2 | 1 | 45 |
| 2.8 | 0 | 22 |



| One | Dist (km) | Day | Commute time (min) |
|-------|-----------|-------|--------------------|
| x_0 | x_1 | x_2 | y |
| 1 | 2.7 | 1 | 25 |
| 1 | 4.1 | 1 | 33 |
| 1 | 1.0 | 0 | 15 |
| 1 | 5.2 | 1 | 45 |
| 1 | 2.8 | 0 | 22 |

Polynomial regression: feature mappings

- Define the feature map:

$$\phi(x) = \begin{pmatrix} 1 \\ x \\ (x)^2 \\ (x)^3 \end{pmatrix}$$

- Polynomial regression model now becomes a linear model w.r.t. the new features

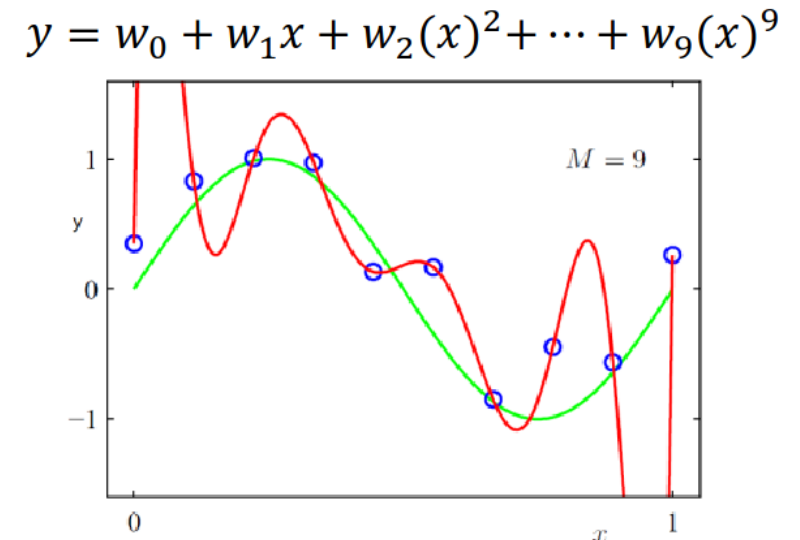
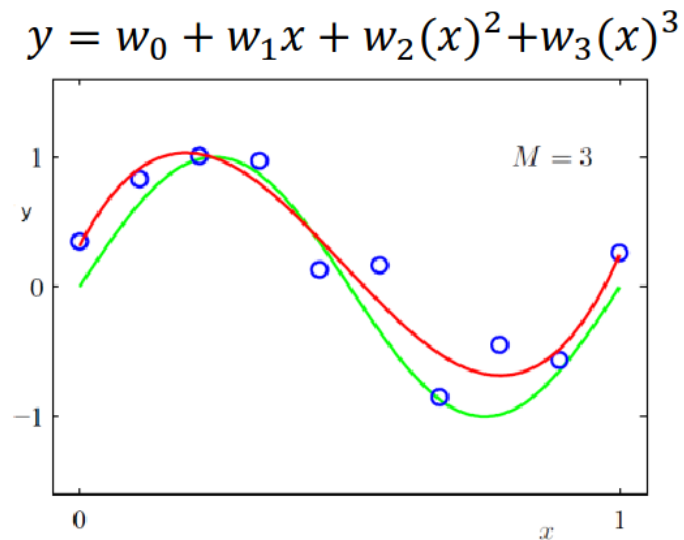
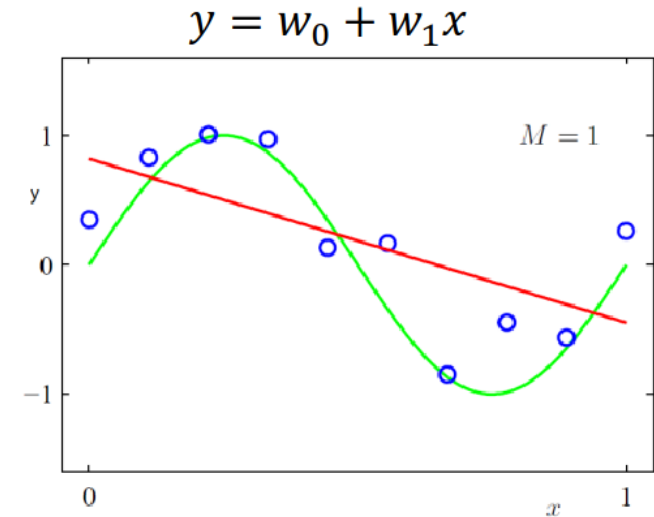
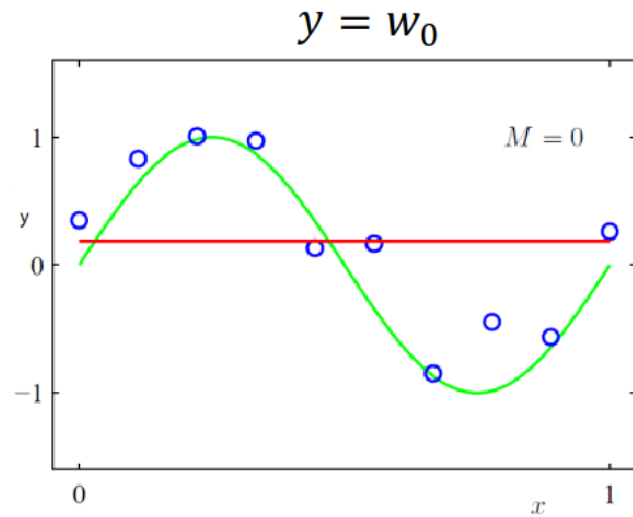
$$f(x) = \mathbf{w}^T \phi(x) = w_0 + w_1 x + w_2 (x)^2 + w_3 (x)^3 = \phi(x)^T \mathbf{w}$$

We've transformed a univariate nonlinear problem to a multivariate linear problem!

- The derivations and algorithms so far in this lecture remain the same!

$$X = \begin{pmatrix} \mathbf{x}^1{}^T \\ \mathbf{x}^2{}^T \\ \vdots \\ \mathbf{x}^n{}^T \end{pmatrix} \mapsto \begin{pmatrix} \phi(x^1)^T \\ \phi(x^2)^T \\ \vdots \\ \phi(x^n)^T \end{pmatrix} = \begin{pmatrix} 1 & x^1 & (x^1)^2 & (x^1)^3 \\ 1 & x^2 & (x^2)^2 & (x^2)^3 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x^n & (x^n)^2 & (x^n)^3 \end{pmatrix} = \bar{X}$$

Polynomial regression: fitting polynomials



Polynomial regression: regularisation

- Regularised least squares regression
 - Given dataset $D = \{(\mathbf{x}^1, y^1), (\mathbf{x}^2, y^2), \dots, (\mathbf{x}^n, y^n)\}$ and a regularisation parameter $\lambda > 0$, find a model to minimise:

$$C(\mathbf{w}) = \underbrace{\frac{1}{2n} (\mathbf{w}^T X^T X \mathbf{w} - 2\mathbf{w}^T X^T \mathbf{y} + \mathbf{y}^T \mathbf{y})}_{\text{fitting to data}} + \underbrace{\frac{\lambda}{2} \|\mathbf{w}\|_2^2}_{\text{regulariser}}$$