

Machine Learning

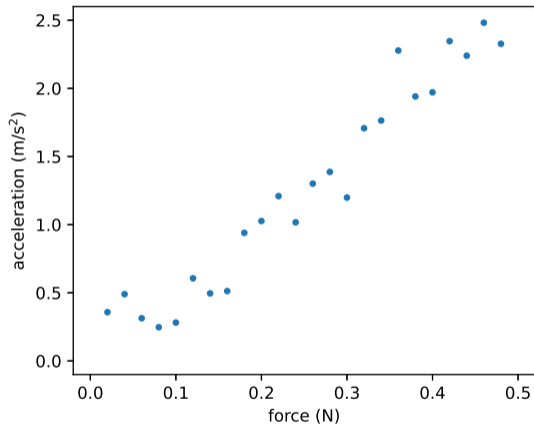
Linear Regression

Hiroshi Shimodaira and Hao Tang

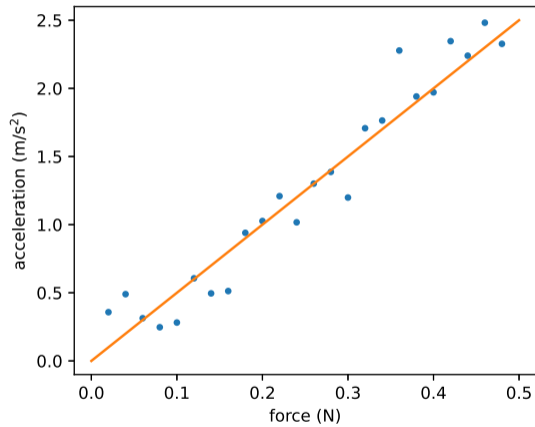
29 January 2024

Ver. 1.1

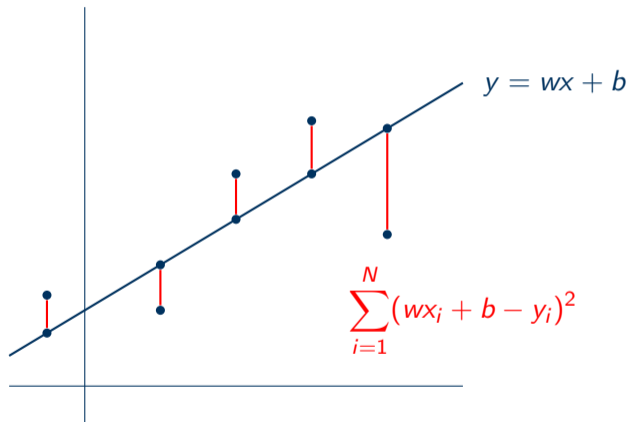
First example



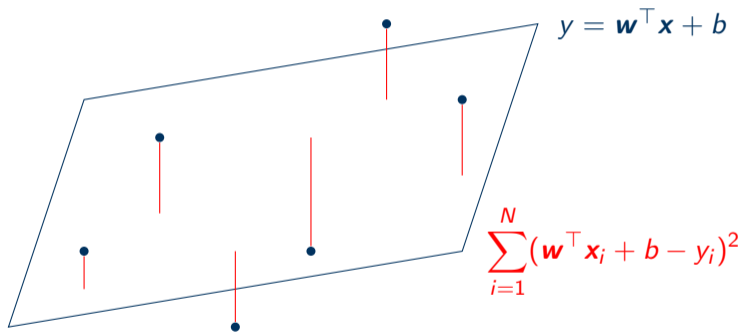
First example



Geometry of linear regression



Geometry of linear regression (cont.)



Linear regression

- $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$: data set
 - $\mathbf{x}_i = [x_{i1} \ \dots \ x_{id}]^\top$: input, features, independent variables
 - $y_i \in \mathbb{R}$: target/dependent variable, ground truth, for \mathbf{x}_i .
- $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$: linear predictor, hyperplane
 - $\mathbf{w} = [w_1 \ \dots \ w_d]^\top$: weights
 - $b \in \mathbb{R}$: bias
 - $\{\mathbf{w}, b\}$: parameters $\dots \ \boldsymbol{\theta} = [b \ \mathbf{w}^\top]^\top$

Linear regression

- Given $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$, find θ such that the mean-squared error (MSE)

$$L = \frac{1}{N} \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i + b - y_i)^2 \quad (1)$$

is minimised.

- The act of finding \mathbf{w} is called training.
- c.f. “least squares” – a parameter estimation method based on MSE or minimising the sum of squares of errors/residuals.

Linear regression: training with MSE

- The goal of linear regression is to solve

$$\min_{\mathbf{w}, b} \frac{1}{N} \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i + b - y_i)^2. \quad (2)$$

- The optimal solution satisfies

$$\frac{\partial L}{\partial b} = 0, \quad \frac{\partial L}{\partial \mathbf{w}} = \left[\frac{\partial L}{\partial w_1} \quad \frac{\partial L}{\partial w_2} \quad \cdots \quad \frac{\partial L}{\partial w_d} \right] = \mathbf{0}. \quad (3)$$

(Is this global optimal? More on this in lectures on optimisation.)

Linear regression: finding the bias b

$$\frac{\partial}{\partial b} \frac{1}{N} \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i + b - y_i)^2 = \frac{2}{N} \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i + b - y_i) \quad (4)$$

$$= 2b - \frac{2}{N} \sum_{i=1}^N (y_i - \mathbf{w}^\top \mathbf{x}_i) = 0 \quad (5)$$

$$b = \frac{1}{N} \sum_{i=1}^N (y_i - \mathbf{w}^\top \mathbf{x}_i) = \frac{1}{N} \sum_{i=1}^N y_i - \mathbf{w}^\top \left(\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \right) = \bar{y} - \mathbf{w}^\top \bar{\mathbf{x}} \quad (6)$$

Linear regression: data centring (mean normalisation)

$$\frac{\partial L}{\partial b} = 0 \implies b = \bar{y} - \mathbf{w}^\top \bar{\mathbf{x}} \quad (7)$$

$$L = \frac{1}{N} \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i + b - y_i)^2 = \frac{1}{N} \sum_{i=1}^N [\mathbf{w}^\top (\mathbf{x}_i - \bar{\mathbf{x}}) - (y_i - \bar{y})]^2 \quad (8)$$

$$= \frac{1}{N} \sum_{i=1}^N (\mathbf{w}^\top \dot{\mathbf{x}}_i - \dot{y}_i)^2 \quad (9)$$

where $\dot{\mathbf{x}}_i = \mathbf{x}_i - \bar{\mathbf{x}}$, $\dot{y}_i = y_i - \bar{y}$

Linear regression: finding the weights w

$$\frac{\partial}{\partial \mathbf{w}} \frac{1}{N} \sum_{i=1}^N (\mathbf{w}^\top \dot{\mathbf{x}}_i - \dot{y}_i)^2 = \frac{2}{N} \sum_{i=1}^N (\mathbf{w}^\top \dot{\mathbf{x}}_i - \dot{y}_i) (\dot{\mathbf{x}}_i) \quad (10)$$

$$= \frac{2}{N} \sum_{i=1}^N ((\mathbf{w}^\top \dot{\mathbf{x}}_i) \dot{\mathbf{x}}_i - \dot{y}_i \dot{\mathbf{x}}_i) \quad (11)$$

Linear regression: finding the weights \mathbf{w} (cont.)

$$\frac{\partial}{\partial \mathbf{w}} \frac{1}{N} \sum_{i=1}^N (\mathbf{w}^\top \dot{\mathbf{x}}_i - \dot{y}_i)^2 = \frac{2}{N} \sum_{i=1}^N ((\mathbf{w}^\top \dot{\mathbf{x}}_i) \dot{\mathbf{x}}_i - \dot{y}_i \dot{\mathbf{x}}_i) \quad (12)$$

$$= \frac{2}{N} \left(\begin{bmatrix} \dot{\mathbf{x}}_1 & \dot{\mathbf{x}}_2 & \cdots & \dot{\mathbf{x}}_N \end{bmatrix} \begin{bmatrix} \mathbf{w}^\top \dot{\mathbf{x}}_1 \\ \mathbf{w}^\top \dot{\mathbf{x}}_2 \\ \vdots \\ \mathbf{w}^\top \dot{\mathbf{x}}_N \end{bmatrix} - \begin{bmatrix} \dot{\mathbf{x}}_1 & \dot{\mathbf{x}}_2 & \cdots & \dot{\mathbf{x}}_N \end{bmatrix} \begin{bmatrix} \dot{y}_1 \\ \dot{y}_2 \\ \vdots \\ \dot{y}_N \end{bmatrix} \right) \quad (13)$$

$$= \frac{2}{N} (\mathbf{X}\mathbf{X}^\top \mathbf{w} - \mathbf{X}\dot{\mathbf{y}}) = \mathbf{0} \quad (14)$$

$$\longrightarrow \mathbf{w} = (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{X}\dot{\mathbf{y}} \quad (15)$$

NB: the definition of \mathbf{X} (which is a $d \times N$ matrix) here is different from the one in the textbook LWLS.

Linear regression - training process

1. Centring

$$\dot{\mathbf{y}} = \begin{bmatrix} y_1 - \bar{y} \\ \vdots \\ y_N - \bar{y} \end{bmatrix}, \quad \mathbf{X} = [\mathbf{x}_1 - \bar{\mathbf{x}} \quad \cdots \quad \mathbf{x}_N - \bar{\mathbf{x}}] \quad (16)$$

2. Computing the weights \mathbf{w} and b

$$\mathbf{w} = (\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{X}\dot{\mathbf{y}} \quad (17)$$

$$b = \bar{y} - \mathbf{w}^\top \bar{\mathbf{x}} \quad (18)$$

NB: $(\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{X}$ is called a Moore-Penrose pseudoinverse of \mathbf{X} .

In practice, we find the solution \mathbf{w} without calculating $(\mathbf{X}\mathbf{X}^\top)^{-1}$

What is \mathbf{XX}^\top ?

- $\mathbf{X} = [\mathbf{x}_1 - \bar{\mathbf{x}}, \dots, \mathbf{x}_N - \bar{\mathbf{x}}]$
- \mathbf{XX}^\top is a $d \times d$ symmetric matrix
- \mathbf{XX}^\top is positive semi-definite, i.e. $\mathbf{x}^\top (\mathbf{XX}^\top) \mathbf{x} \geq 0$ for any $\mathbf{x} \in \mathbb{R}^d$
NB: Eigen values of a positive semi-definite matrix are non-negative, i.e. $\lambda_i \geq 0$ for $i = 1, \dots, d$
- $\mathbf{C} = \frac{1}{N} \mathbf{XX}^\top$ is called a **covariance matrix**
 - $\mathbf{C} = (\sigma_{ij})$: σ_{ii} is the (population) variance of i -th dimension of data, σ_{ij} is the covariance between i -th and j -th dimensions of data.
 - used in many areas, e.g. multivariate normal distributions, principal component analysis (PCA)
- $\det(\mathbf{C}) = \prod_{i=1}^d \lambda_i$ and $\text{tr}(\mathbf{C}) = \sum_{i=1}^d \lambda_i$, where λ_i is the i -th eigen value of \mathbf{C}
- $\det(\mathbf{C}) = 0$ and $\text{rk}(\mathbf{C}) < d$ if $N \leq d$

Features

$$y = \mathbf{w}^\top \mathbf{x} + b = [\mathbf{w}^\top \quad b] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix} = \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}^\top \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix} = \mathbf{w}'^\top \mathbf{x}' \quad (19)$$

- Fitting $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$ is equivalent to appending 1 to \mathbf{x} and fitting $f(\mathbf{x}') = \mathbf{w}'^\top \mathbf{x}'$.
- The 1 can be seen as a feature independent of the input.

Features

- Suppose we have a data point $\mathbf{x} = [x_1 \ x_2 \ x_3]^\top$.
- The data point after appending 1 becomes

$$[1 \ x_1 \ x_2 \ x_3]^\top \quad (20)$$

- The data point after appending 1 and quadratic terms becomes

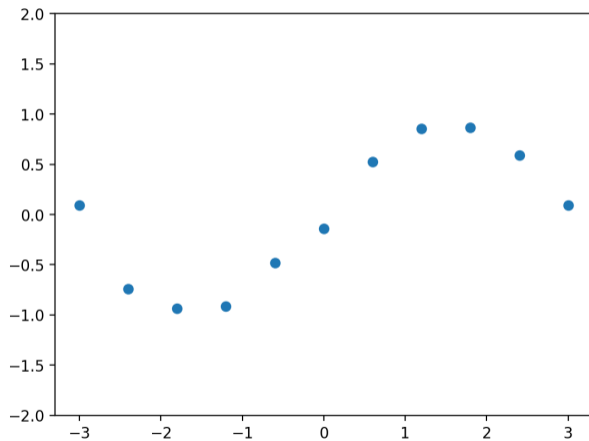
$$\phi(\mathbf{x}) = [1 \ x_1 \ x_2 \ x_3 \ x_1x_2 \ x_2x_3 \ x_1x_3 \ x_1^2 \ x_2^2 \ x_3^2]^\top \quad (21)$$

- The function $f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$ is a polynomial.

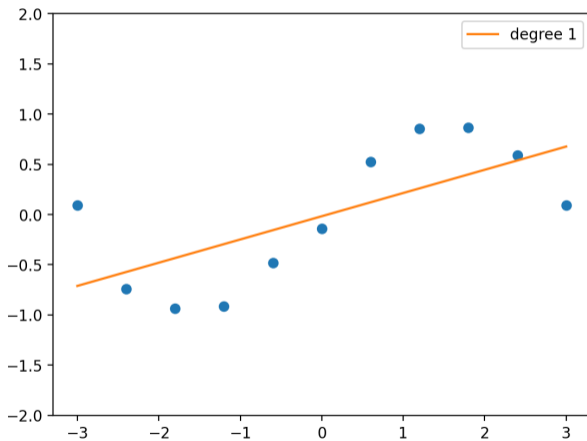
Linear regression with feature transformation

- We call ϕ a feature function.
- In general, ϕ can be any function.
- Instead of $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$, we now have $f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$.
- Instead of $\mathbf{X} = [\mathbf{x}_1 \quad \mathbf{x}_2 \quad \cdots \quad \mathbf{x}_N]$, we have $\Phi = [\phi(\mathbf{x}_1) \quad \phi(\mathbf{x}_2) \quad \cdots \quad \phi(\mathbf{x}_N)]$
- The optimal solution for linear regression becomes $\mathbf{w} = (\Phi\Phi^\top)^{-1}\Phi\mathbf{y}$.

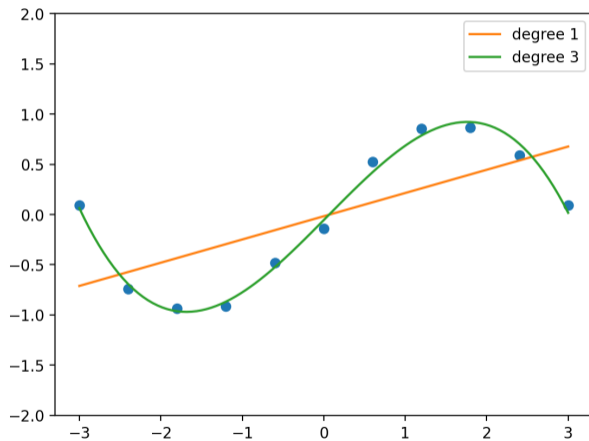
Examples



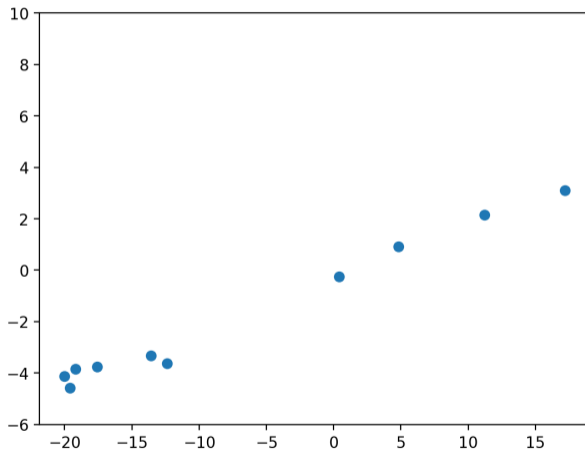
Examples



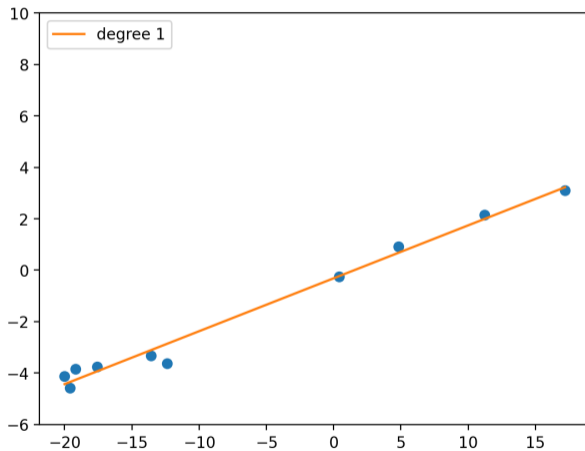
Examples



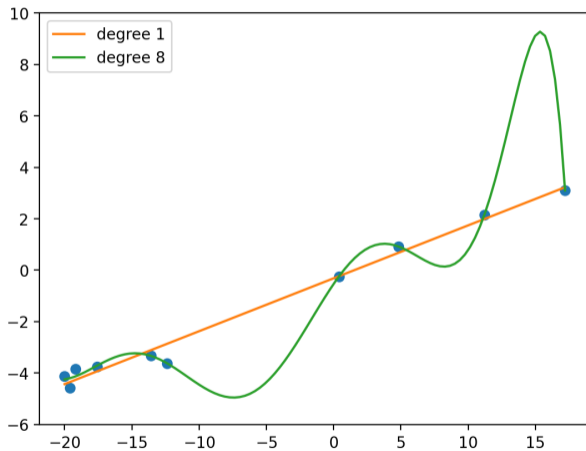
Examples



Examples



Examples



Linear regression

- A “linear” regression model is linear in the parameters \mathbf{w} (i.e. linear combination between the parameters and features), **not** the features.
- A linear regression model can fit an arbitrary nonlinear function.
- What are the “right” features?
- What does it mean for the program $\mathbf{w}^\top \phi(\mathbf{x})$ we write with data to be “correct”? (Is it right to use a complex nonlinear transformation $\phi(\mathbf{x})$?)

A probabilistic interpretation

- Assume we cannot get a perfect fit because of noise.
- In particular, we assume the noise is additive and Gaussian.
- In other words, $y_i = \mathbf{w}^\top \phi(\mathbf{x}_i) + \epsilon_i$, where $\epsilon_i \sim \mathcal{N}(0, 1)$.
- If $\epsilon_i \sim \mathcal{N}(0, 1)$, then $y_i \sim \mathcal{N}(\mathbf{w}^\top \phi(\mathbf{x}_i), 1)$.
- The log-likelihood of \mathbf{w} is

$$\log \prod_{i=1}^N \frac{1}{\sqrt{2\pi}} \exp \left(-\frac{1}{2} (y_i - \mathbf{w}^\top \phi(\mathbf{x}_i))^2 \right) \quad (22)$$

A probabilistic interpretation

- Log-likelihood of \mathbf{w}

$$\sum_{i=1}^N \left[-\frac{1}{2} \log(2\pi) - \frac{1}{2} (y_i - \mathbf{w}^\top \phi(\mathbf{x}_i))^2 \right] \quad (23)$$

- Mean-squared error

$$\frac{1}{N} \sum_{i=1}^N (y_i - \mathbf{w}^\top \phi(\mathbf{x}_i))^2 \quad (24)$$

- The maximum likelihood estimator is the optimal solution for MSE.

Practical issues

- The complexity of computing $(\Phi\Phi^\top)\Phi\mathbf{y}$ is $O(N^3)$, where N is the number of samples.
- The runtime is not particularly suitable for large data sets.
- Instead of solving $\min_{\mathbf{w}} L$ exactly, could we find an approximate solution?
- In exchange, could we have an algorithm that scales better than $O(N^3)$?
- Not all problems have closed-form solutions for $\frac{\partial L}{\partial \mathbf{w}}$ anyways.
- What if there are outliers?

Linear regression

- We write a program $f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$ with $\mathbf{w} = (\Phi\Phi^\top)^{-1}\Phi\mathbf{y}$.
- In what sense is this program “correct”?

Linear regression using matrix calculus

- The mean-squared error can be written compactly as

$$L = \|\Phi^T \mathbf{w} - \mathbf{y}\|_2^2. \quad (25)$$

- We can expand the mean-squared error as

$$L = \|\Phi^T \mathbf{w} - \mathbf{y}\|_2^2 = (\Phi^T \mathbf{w} - \mathbf{y})^T (\Phi^T \mathbf{w} - \mathbf{y}) = \mathbf{w}^T \Phi \Phi^T \mathbf{w} - 2\mathbf{y}^T \Phi^T \mathbf{w} + \mathbf{y}^T \mathbf{y}. \quad (26)$$

- Solving the optimal solution gives

$$\frac{\partial L}{\partial \mathbf{w}} = (\Phi \Phi^T + (\Phi \Phi^T)^T) \mathbf{w} - 2\Phi \mathbf{y} = \mathbf{0} \implies \mathbf{w} = (\Phi \Phi^T)^{-1} \Phi \mathbf{y}. \quad (27)$$

Topics not covered

- Choices of features \mathbf{x} (feature selection)
- Interpretations of the model parameters θ
- Collinearity
- Heteroscedasticity
- Other linear regression models (e.g. ridge regression, LASSO, Bayesian linear regression)
- Multiple linear regression
- Relationships with neural networks
- Relationships with principal component analysis (PCA)

Quizzes

1. What is the number of dimensions of the hyperplane formed by linear regression?
2. Give detailed derivations for Eqs. (13) and (14).
3. Show that $\mathbf{X}\mathbf{X}^\top$ is positive semi-definite.
4. Using $\mathbf{x}' = (1, \mathbf{x}^\top)^\top$ instead of \mathbf{x} , rewrite Eqs. (1), \dots , (15).