

Inductive Learning

Readings: ( :chapt. 3)

1

Pure inductive learning

- Let f be an *unknown* target function

Problem: find a hypothesis h
 such that $h \approx f$
 given a training set of examples

An *example* is a pair $(x, f(x))$

2

Pure inductive learning

- The previous problem is an *ill-defined problem*

In general, data are not enough to identify a unique hypothesis.

19-Sep-24

<http://tiny.cc/ml2024-25>

3

3

Pure inductive learning

- A binary function of two binary variables (x_1, x_2)

x_1	x_2	h_1	h_2	h_3	h_4	h_5	h_6	h_7	h_8	h_9	h_{10}	h_{11}	h_{12}	h_{13}	h_{14}	h_{15}	h_{16}
0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1
0	1	0	0	0	0	1	1	1	1	0	0	0	0	1	1	1	1
1	0	0	0	1	1	0	0	1	1	0	0	1	1	0	0	1	1
1	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1

19-Sep-24

<http://tiny.cc/ml2024-25>

12 - 4

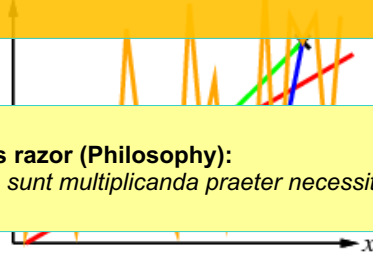
4

Pure inductive learning

- How to choose among different hypothesis (models)?

As the problem is ill-defined, an *inductive bias*, i.e., a set of à priori assumptions is required.

Ockham's razor (Philosophy):
"entia non sunt multiplicanda praeter necessitatem"



19-Sep-24

<http://tiny.cc/ml2024-25>

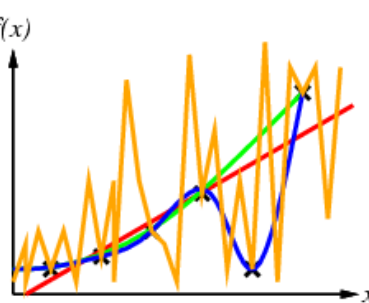
12 - 5

5

Ockham's razor – Why?

■ Advantages

- There are less simple hypothesis than complex ones.
- highly probable that a sufficiently complex hypothesis will fit the data
- An simple hypothesis that fits data is less probable to be a statistical coincidence.



19-Sep-24

<http://tiny.cc/ml2024-25>

12 - 6

6

Inductive learning hypothesis

- **Inductive learning hypothesis:** Any hypothesis found to approximate the target function well over a sufficient large set of training examples will also approximate the target function well over other unobserved examples.

19-Sep-24

<http://tiny.cc/ml2024-25>

12 – 7

7

Inductive Learning

- **Generalization:** What's the model's performance on new data?
- **Overfitting:** \mathcal{H} is more complex than f
- **Underfitting:** \mathcal{H} is less complex than f

19-Sep-24

<http://tiny.cc/ml2024-25>

8

8

To take away

- Inductive learning
 - Pure vs. biased inductive learning
 - Ockham's razor
 - Inductive Learning hypothesis

19-Sep-24

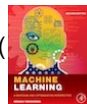
<http://tiny.cc/ml2024-25>

12 – 9

9

Learning in parametric modeling (Simple linear regression)

Readings: (

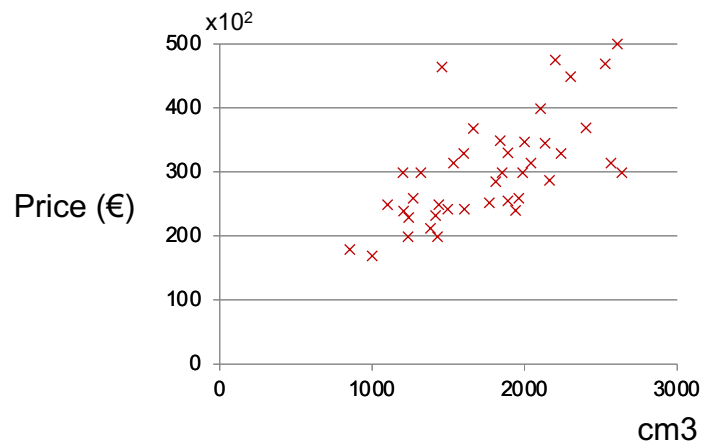


:chapt. 3)

10

Working example:

Car's price as a function of engine volume



19-Sep-24

<http://tiny.cc/ml2024-25>

12 – 11

11

Some notation

Cm3 (x)	Price in € (y)
2104	460
852	178
...	...

$D = \{(x^{(i)}, y^{(i)}) \mid i=1, \dots, m\}$ – Data set

m - Number of examples (elements) in the data set

$x^{(i)}$ - input variable (feature) of the i -th example

$y^{(i)}$ – desired output variable of the i -th example

19-Sep-24

<http://tiny.cc/ml2024-25>

12 – 12

12

Hypothesis

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

where θ_0, θ_1 are real-valued coefficients or weights

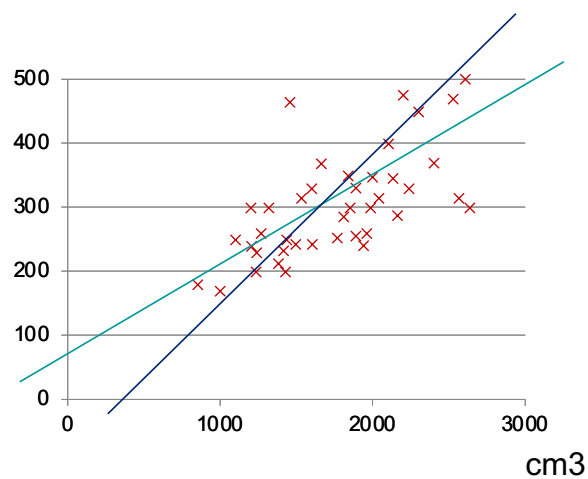
19-Sep-24

<http://tiny.cc/ml2024-25>

12 – 13

13

Working example:
Car's price as a function of its engine



14

Performance function (or loss function)

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m \underbrace{(h_{\theta}(x^{(i)}) - y^{(i)})}_{\text{residual}}^2$$

Objective: minimize $J(\theta_0, \theta_1)$
 θ_0, θ_1

19-Sep-24

<http://tiny.cc/ml2024-25>

12 – 15

15

Y-intercept and sample slope

- Hypothesis: $h_{\theta}(x) = \theta_0 + \theta_1 x$
- Y-Intercept: $\theta_0 = \bar{y} - \theta_1 \bar{x}$
- Sample' slope $\theta_1 = \frac{SS_{xy}}{SS_{xx}} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$

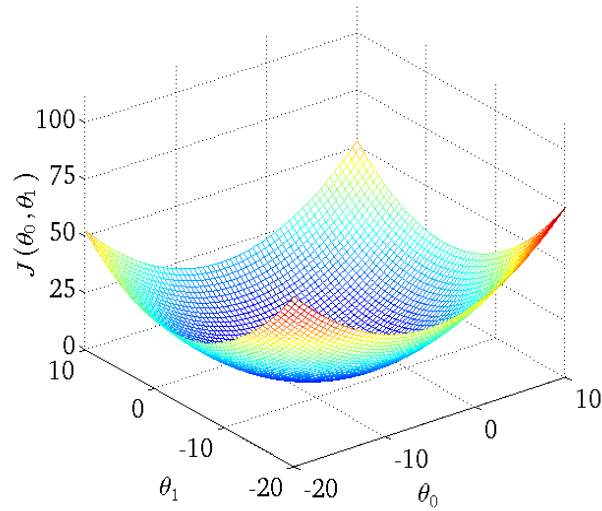
19-Sep-24

<http://tiny.cc/ml2024-25>

12 – 16

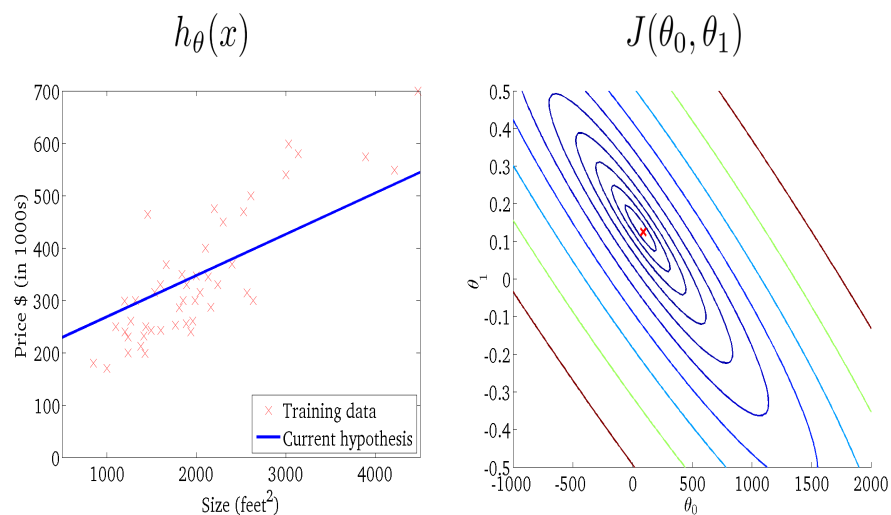
16

J as function of θ_0 and θ_1



17

Contour of $J(\theta_0, \theta_1)$



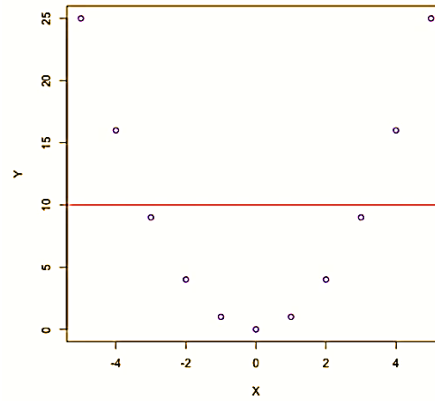
19-Sep-24

<http://tiny.cc/ml2024-25>

12 – 18

18

Linear regression of a quadratic function



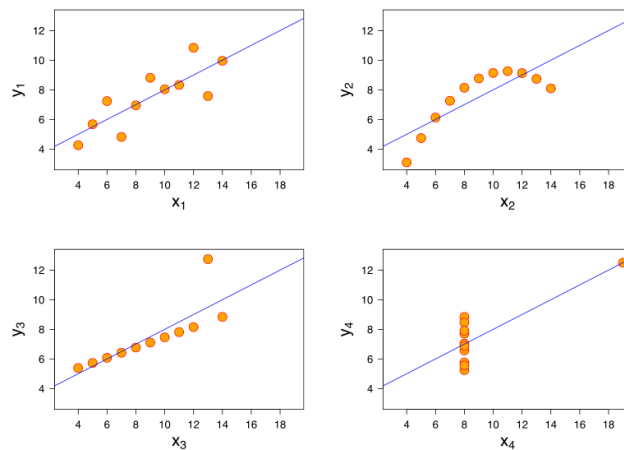
19-Sep-24

<http://tiny.cc/ml2024-25>

12 – 19

19

F. J. Anscombe's quartet (1973)



19-Sep-24

<http://tiny.cc/ml2024-25>

12 – 20

20

Residuals

- Residuals are the difference between the observed (y) and predicted responses (h).
- For the normal error regression model, we assume that the error term is normally distributed with null mean and constant variance.
- If the model is appropriate for the data, this should be reflected in the residuals. In particular, the residual vs. input graphic should be centered in zero and show no tendency.

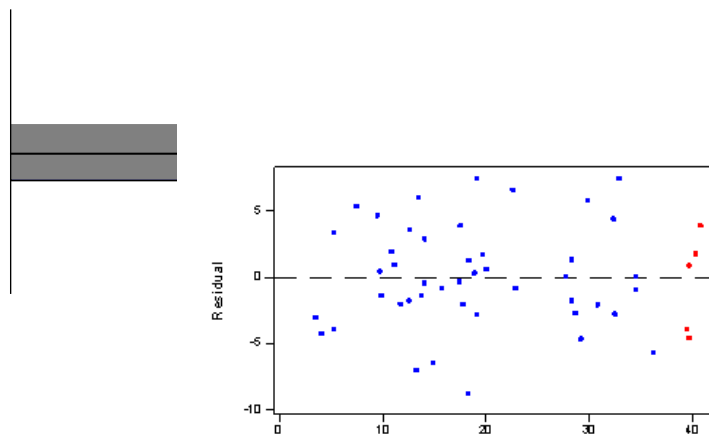
19-Sep-24

<http://tiny.cc/ml2024-25>

12 – 21

21

Residuals, the pattern we would like to see



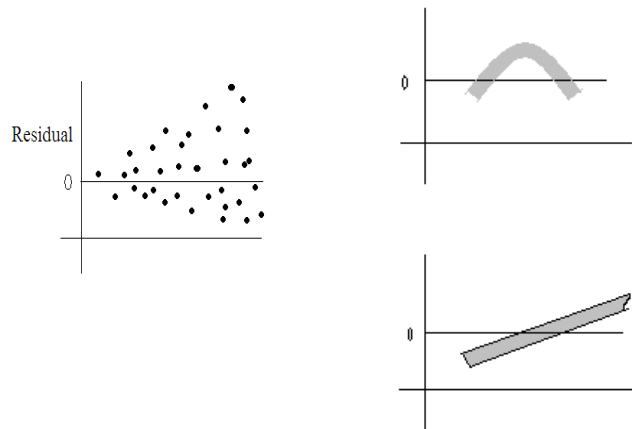
19-Sep-24

<http://tiny.cc/ml2024-25>

12 – 22

22

Residuals with some kind of anomalies



19-Sep-24

<http://tiny.cc/ml2024-25>

23

23

To take away

- Inductive learning
 - Pure vs. biased inductive learning
 - Ockham's razor
 - Inductive Learning hypothesis
- Simple Linear regression
 - Notation
 - Performance or total loss function
 - Y-Intercept, slope, and residuals
 - Limitations

19-Sep-24

<http://tiny.cc/ml2024-25>

12 – 24

24

Recalling the gradient method

25

Problem formulation

- Find the x^* such that minimizes

$$f : R^n \rightarrow R$$

where f is a smooth C^1 function

19-Sep-24

<http://tiny.cc/ml2024-25>

16 – 26

26

The gradient

$$f : \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\nabla f(x_1, \dots, x_n) := \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right)^T$$

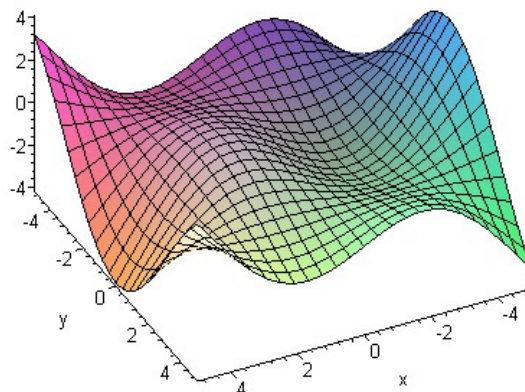
19-Sep-24

<http://tiny.cc/ml2024-25>

16 – 27

27

Application example



$$f := (x, y) \rightarrow \cos\left(\frac{1}{2}x\right) \cos\left(\frac{1}{2}y\right) x$$

19-Sep-24

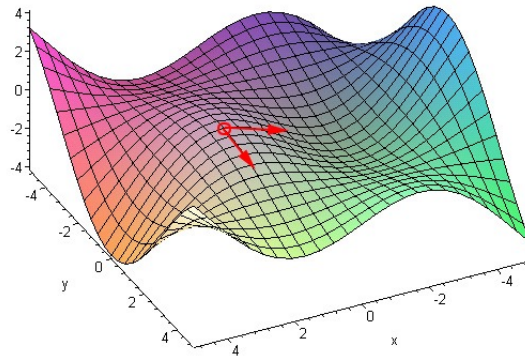
12 – 28

28

Directional derivatives

$$\frac{\partial f(x, y)}{\partial y}$$

$$\frac{\partial f(x, y)}{\partial x}$$



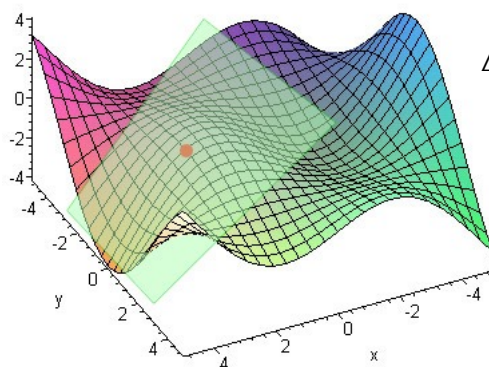
19-Sep-24

12 - 29

29

The gradient properties

- The gradient defines a (hyper) plane approximating the function infinitesimally



$$\Delta z = \frac{\partial f}{\partial x} \cdot \Delta x + \frac{\partial f}{\partial y} \cdot \Delta y$$

19-Sep-24

12 - 30

30

The gradient properties

- **Proposition:** let $f: R^n \rightarrow R$ be a smooth C^1 function around p , if f has local minimum (or maximum) at p then,

$$(\nabla f)_p = \bar{0}$$

19-Sep-24

<http://tiny.cc/ml2024-25>

16 - 31

31

Gradient (steepest descent) algorithm

$i=0$

$x_0 \in R^n$

Repeat until $\nabla f(x_i) = 0$

1. compute *search direction* $h_i = -\nabla f(x_i)$

2. update $x_{i+1} = x_i + \eta \cdot h_i$

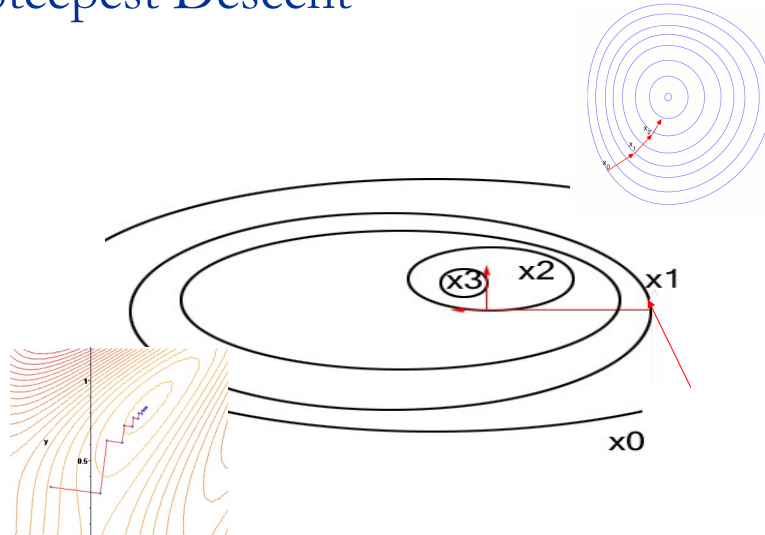
19-Sep-24

<http://tiny.cc/ml2024-25>

16 - 32

32

Steepest Descent



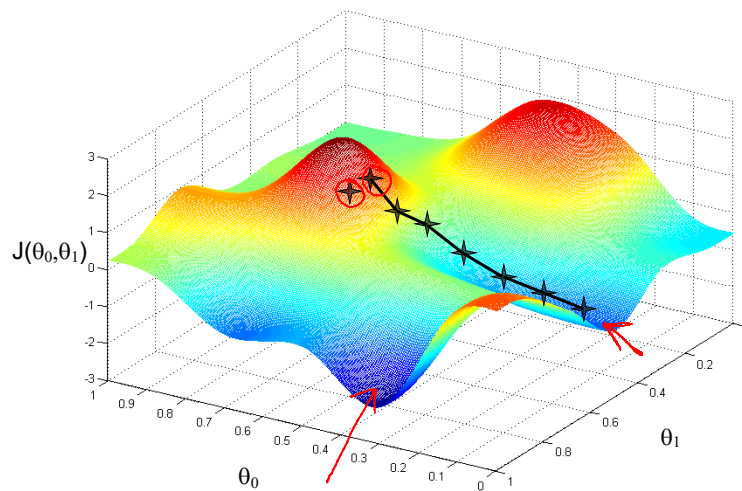
19-Sep-24

<http://tiny.cc/ml2024-25>

16 – 33

33

The effect of initial conditions



19-Sep-24

<http://tiny.cc/ml2024-25>

16 – 34

34

Local minima and weight initialization

- The steepest decent gradient finds a minimum, not necessarily a global minimum,
- Run the algorithm N times with small different random initial values for the weights.

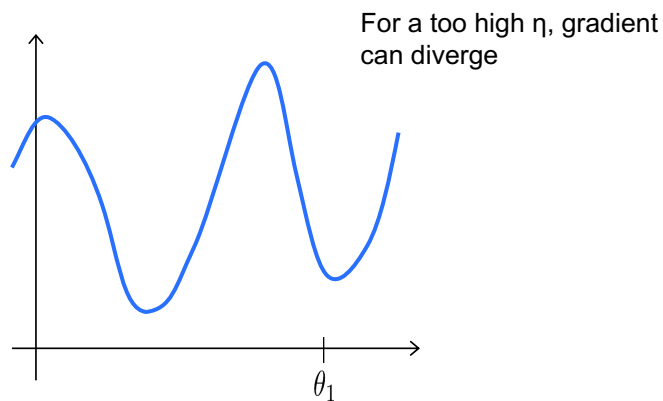
19-Sep-24

<http://tiny.cc/ml2024-25>

18 – 35

35

Learning step size parameter, η



19-Sep-24

<http://tiny.cc/ml2024-25>

16 – 36

36

Gradient and regression in brief

Gradient

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

 (for $j = 1$ and $j = 0$)
 }

Linear regression model

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

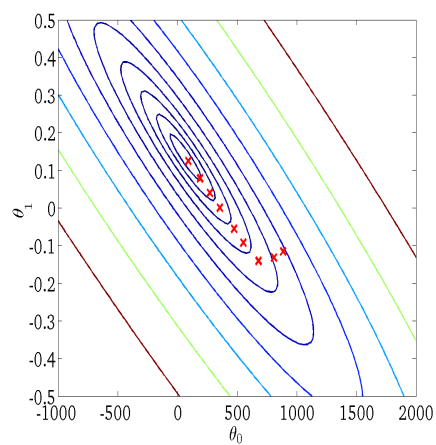
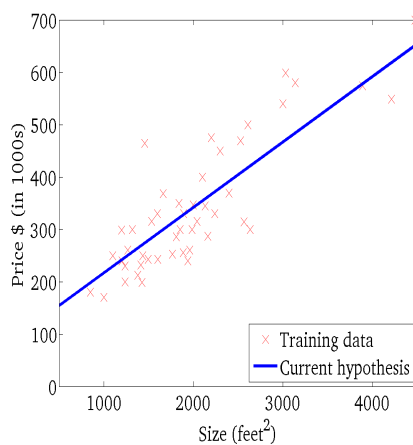
19-Sep-24

<http://tiny.cc/ml2024-25>

16 - 37

37

Evolution of $J(\theta_0, \theta_1)$

<http://tiny.cc/ml2024-25>

38

To take away today