

# Linear regression

Petr Pošík

<b>Linear regression</b>	2
Regression	3
Notation remarks	4
Train, apply	5
1D regression	6
LSM	7
Minimizing $J(w, T)$	8
Multivariate linear regression	9

**Linear regression**

**Regression task** is a supervised learning task, i.e.

- a training (multi)set  $T = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(|T|)}, y^{(|T|)})\}$  is available, where
- the labels  $y^{(i)}$  are *quantitative*, often continuous (as opposed to classification tasks where  $y^{(i)}$  are nominal).
- Its purpose is to model the relationship between independent variables (inputs)  $\mathbf{x} = (x_1, \dots, x_D)$  and the dependent variable (output)  $y$ .

**Linear regression** is a particular regression model which assumes (and learns) linear relationship between the inputs and the output:

$$\hat{y} = h(\mathbf{x}) = w_0 + w_1x_1 + \dots + w_Dx_D = w_0 + \langle \mathbf{w}, \mathbf{x} \rangle = w_0 + \mathbf{x}\mathbf{w}^T,$$

where

- $\hat{y}$  is the model *prediction (estimate)* of the true value  $y$ ,
- $h(\mathbf{x})$  is the linear model (a *hypothesis*),
- $w_0, \dots, w_D$  are the coefficients of the linear function,  $w_0$  is the *bias*, organized in a row vector  $\mathbf{w}$ ,
- $\langle \mathbf{w}, \mathbf{x} \rangle$  is a *dot product* of vectors  $\mathbf{w}$  and  $\mathbf{x}$  (scalar product),
- which can be also computed as a matrix product  $\mathbf{x}\mathbf{w}^T$  if  $\mathbf{w}$  and  $\mathbf{x}$  are row vectors.

**Notation remarks**

**Homogeneous coordinates:** If we add "1" as the first element of  $\mathbf{x}$  so that  $\mathbf{x} = (1, x_1, \dots, x_D)$ , then we can write the linear model in an even simpler form (without the explicit bias term):

$$\hat{y} = h(\mathbf{x}) = w_0 \cdot 1 + w_1x_1 + \dots + w_Dx_D = \langle \mathbf{w}, \mathbf{x} \rangle = \mathbf{x}\mathbf{w}^T.$$

**Matrix notation:** If we organize the data into matrix  $\mathbf{X}$  and vector  $\mathbf{y}$ , such that

$$\mathbf{X} = \begin{pmatrix} 1 & \mathbf{x}^{(1)} \\ \vdots & \vdots \\ 1 & \mathbf{x}^{(|T|)} \end{pmatrix} \quad \text{and} \quad \mathbf{y} = \begin{pmatrix} y^{(1)} \\ \vdots \\ y^{(|T|)} \end{pmatrix},$$

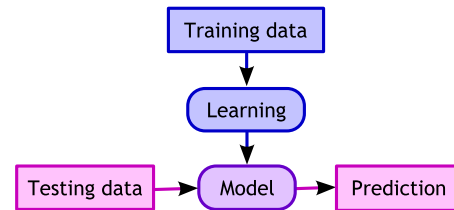
and similarly with  $\hat{\mathbf{y}}$ , then we can write a batch computation of predictions for all data in  $\mathbf{X}$  as

$$\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}^T.$$

## Two operation modes

Any ML model has 2 operation modes:

1. learning (training, fitting) and
2. application (testing, making predictions).



The model  $h$  can be viewed as a function of 2 variables:  $h(x, w)$ .

**Model application:** If the model is given ( $w$  is fixed), we can manipulate  $x$  to make predictions:

$$\hat{y} = h(x, w) = h_w(x).$$

**Model learning:** If the data is given ( $T$  is fixed), we can manipulate the model parameters  $w$  to fit the model to the data:

$$w^* = \underset{w}{\operatorname{argmin}} J(w, T).$$

**How to train the model?**

## Simple (univariate) linear regression

**Simple (univariate) regression** deals with cases where  $x^{(i)} = x^{(i)}$ , i.e. the examples are described by a single feature (they are 1-dimensional).

**Fitting a line to data:**

- find parameters  $w_0, w_1$  of a linear model  $\hat{y} = w_0 + w_1x$
- given a training (multi)set  $T = \{(x^{(i)}, y^{(i)})\}_{i=1}^{|T|}$ .

How to fit depending on the number of training examples:

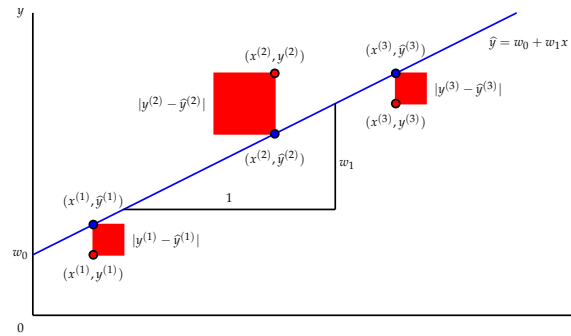
- Given a single example (1 equation, 2 parameters)  
⇒ infinitely many linear function can be fitted.
- Given 2 examples (2 equations, 2 parameters)  
⇒ exactly 1 linear function can be fitted.
- Given 3 or more examples ( $> 2$  equations, 2 parameters)  
⇒ no line can be fitted without any error  
⇒ a line which minimizes the “size” of error  $y - \hat{y}$  can be fitted:

$$w^* = (w_0, w_1) = \underset{w_0, w_1}{\operatorname{argmin}} J(w_0, w_1, T).$$

## The least squares method

The **least squares method (LSM)** suggests to choose such parameters  $w$  which minimize the *mean squared error*

$$J(w) = \frac{1}{|T|} \sum_{i=1}^{|T|} (y^{(i)} - \hat{y}^{(i)})^2 = \frac{1}{|T|} \sum_{i=1}^{|T|} (y^{(i)} - h_w(x^{(i)}))^2.$$

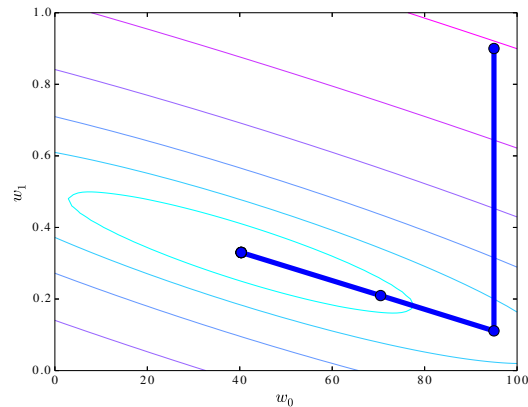
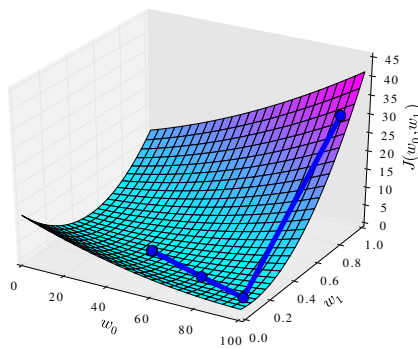


**Explicit solution:**

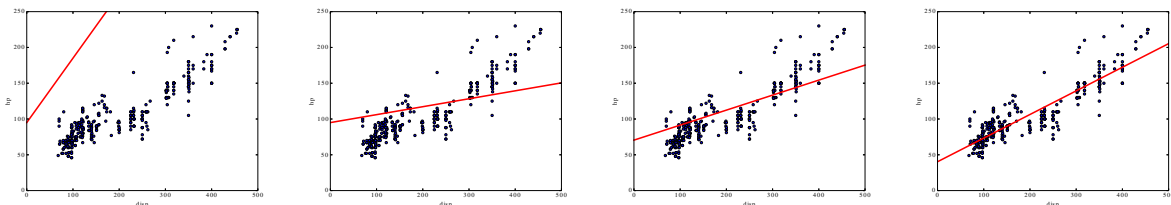
$$w_1 = \frac{\sum_{i=1}^{|T|} (x^{(i)} - \bar{x})(y^{(i)} - \bar{y})}{\sum_{i=1}^{|T|} (x^{(i)} - \bar{x})^2} = \frac{s_{xy}}{s_x^2} \quad w_0 = \bar{y} - w_1 \bar{x}$$

## Universal fitting method: minimization of cost function $J$

The landscape of  $J$  in the space of  $w_0$  and  $w_1$ :



Gradually better linear models found by an optimization method (BFGS):



## Multivariate linear regression

**Multivariate linear regression** deals with cases where  $\mathbf{x}^{(i)} = (x_1^{(i)}, \dots, x_D^{(i)})$ , i.e. the examples are described by more than 1 feature (they are  $D$ -dimensional).

### Model fitting:

- find parameters  $\mathbf{w} = (w_1, \dots, w_D)$  of a linear model  $\hat{y} = \mathbf{x}\mathbf{w}^T$
- given the training (multi)set  $T = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^{|T|}$ .
- The model is a *hyperplane* in the  $D + 1$ -dimensional space.

Fitting methods:

1. Numeric optimization of  $J(\mathbf{w}, T)$ :
  - Works as for simple regression, it only searches a space with more dimensions.
  - Sometimes one need to tune some parameters of the optimization algorithm to work properly (learning rate in gradient descent, etc.).
  - May be slow (many iterations needed), but works even for very large  $D$ .
2. **Normal equation:**

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

- Method to solve for the optimal  $\mathbf{w}^*$  analytically!
- No need to choose optimization algorithm parameters.
- No iterations.
- Needs to compute  $(\mathbf{X}^T \mathbf{X})^{-1}$ , which is  $O(D^3)$ . Slow, or intractable, for large  $D$ .