

Linear Regression II, SGD

Milan Straka

 October 10, 2022



Charles University in Prague
Faculty of Mathematics and Physics
Institute of Formal and Applied Linguistics



unless otherwise stated

Linear Regression

Given an input value $\mathbf{x} \in \mathbb{R}^D$, **linear regression** computes predictions as:

$$y(\mathbf{x}; \mathbf{w}, b) = x_1 w_1 + x_2 w_2 + \dots + x_D w_D + b = \sum_{i=1}^D x_i w_i + b = \mathbf{x}^T \mathbf{w} + b.$$

The *bias* b can be considered one of the *weights* \mathbf{w} if convenient.

We train the weights by minimizing an **error function** between the real target values and their predictions, notably *sum of squares*:

$$\frac{1}{2} \sum_{i=1}^N (y(\mathbf{x}_i; \mathbf{w}) - t_i)^2$$

There are various approaches to minimize it, but for linear regression an explicit solution exists:

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}.$$

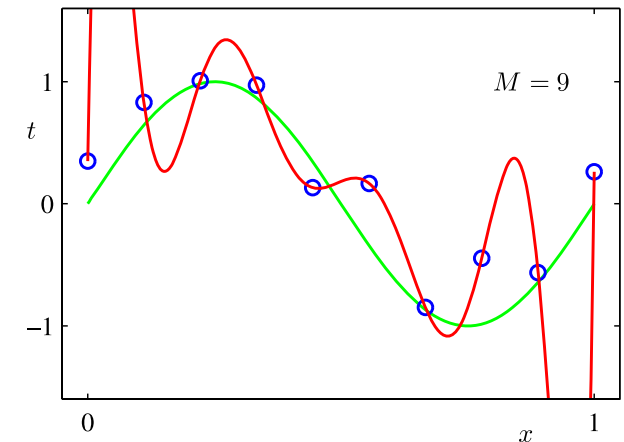
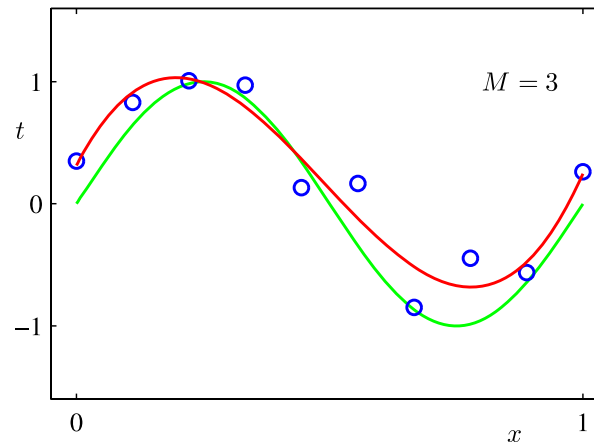
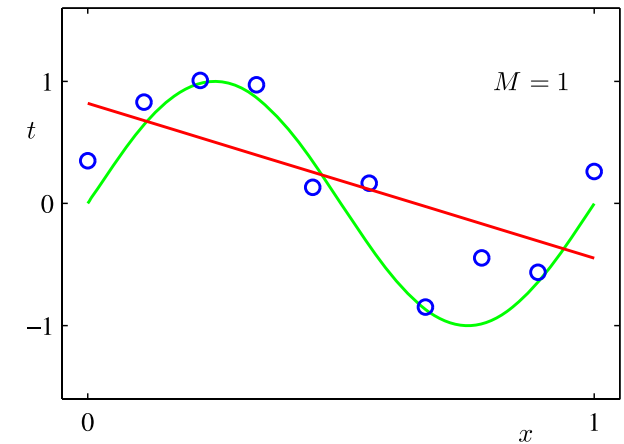
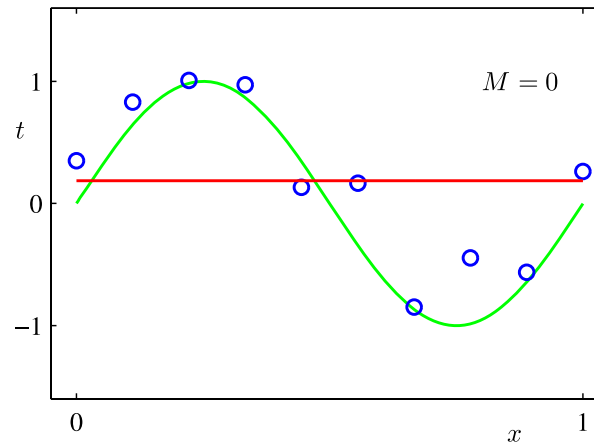
Linear Regression Example

Assume we want to predict a $t \in \mathbb{R}$ for a given $x \in \mathbb{R}$. If we train the linear regression with “raw” input vectors $\mathbf{x} = (x)$, only straight lines could be modeled.

However, if we consider input vectors $\mathbf{x} = (x^0, x^1, \dots, x^M)$ for a given $M \geq 0$, the linear regression is able to model polynomials of degree M , because the prediction is then computed as

$$w_0x^0 + w_1x^1 + \dots + w_Mx^M.$$

Therefore, the weights are the coefficients of a polynomial of degree M .



Linear Regression Example

To plot the error, the *root mean squared error* $\text{RMSE} = \sqrt{\text{MSE}}$ is frequently used.

The displayed error nicely illustrates two main challenges in machine learning:

- *underfitting*
- *overfitting*

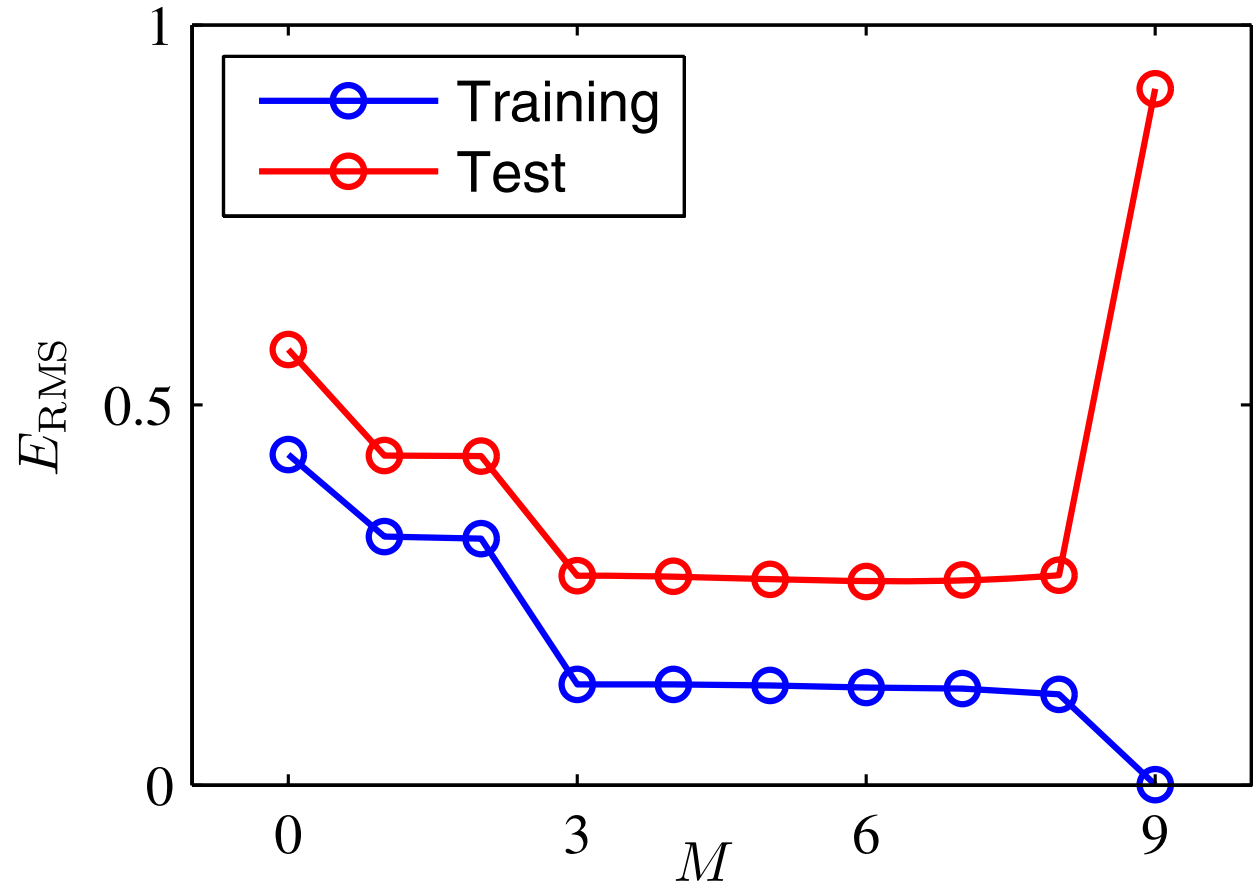


Figure 1.5 of Pattern Recognition and Machine Learning.

Model Capacity

We can control whether a model underfits or overfits by modifying its **capacity**.

- representational capacity
- effective capacity

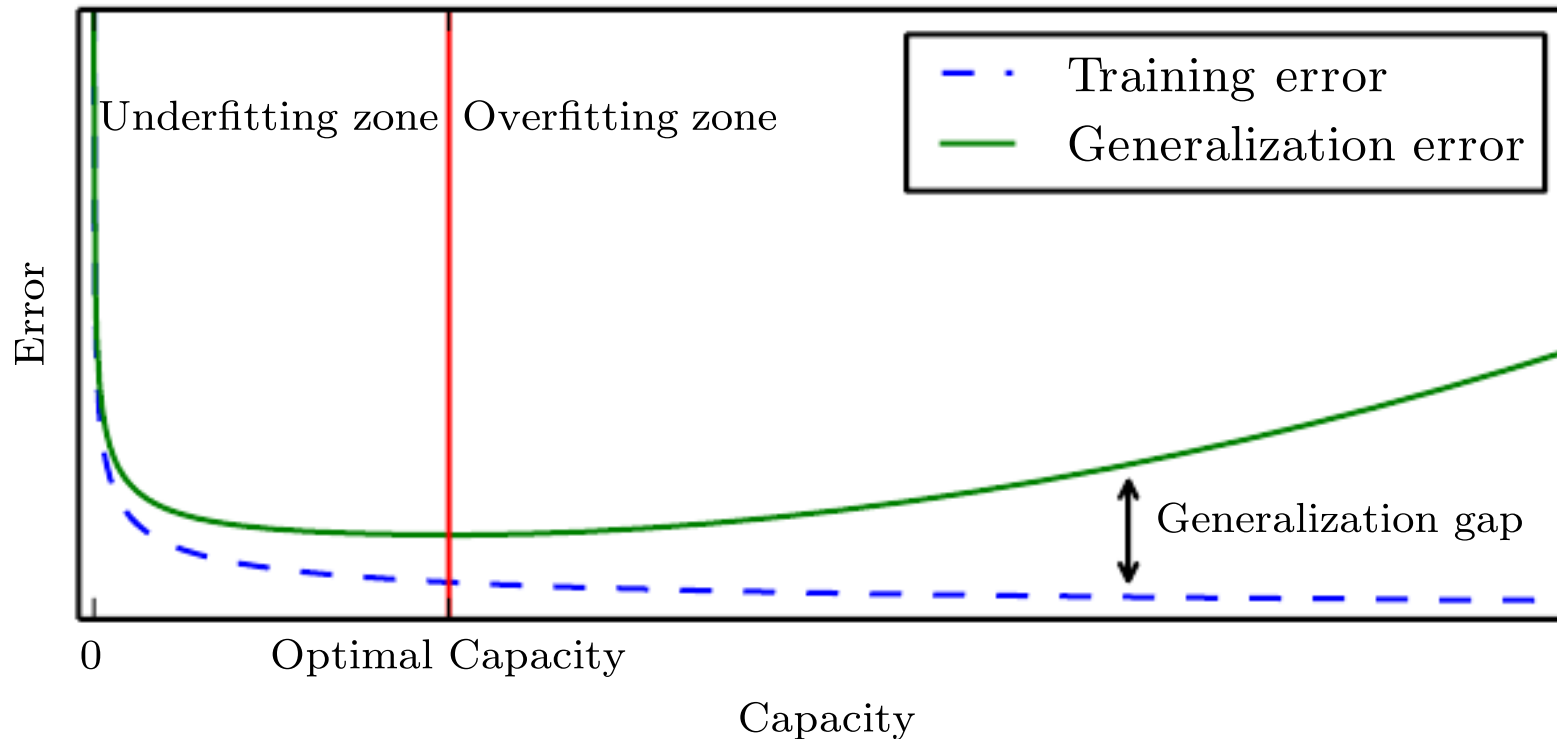


Figure 5.3 of "Deep Learning" book, <https://www.deeplearningbook.org>

Linear Regression Overfitting

Note that employing more data usually alleviates overfitting (the relative capacity of the model is decreased).

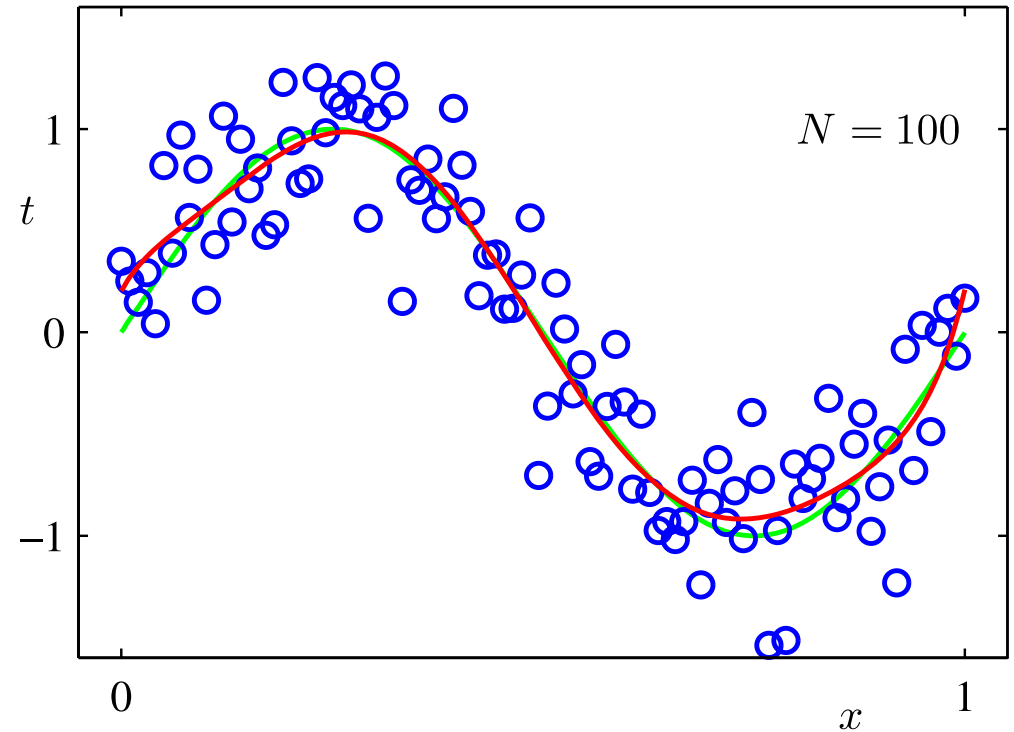
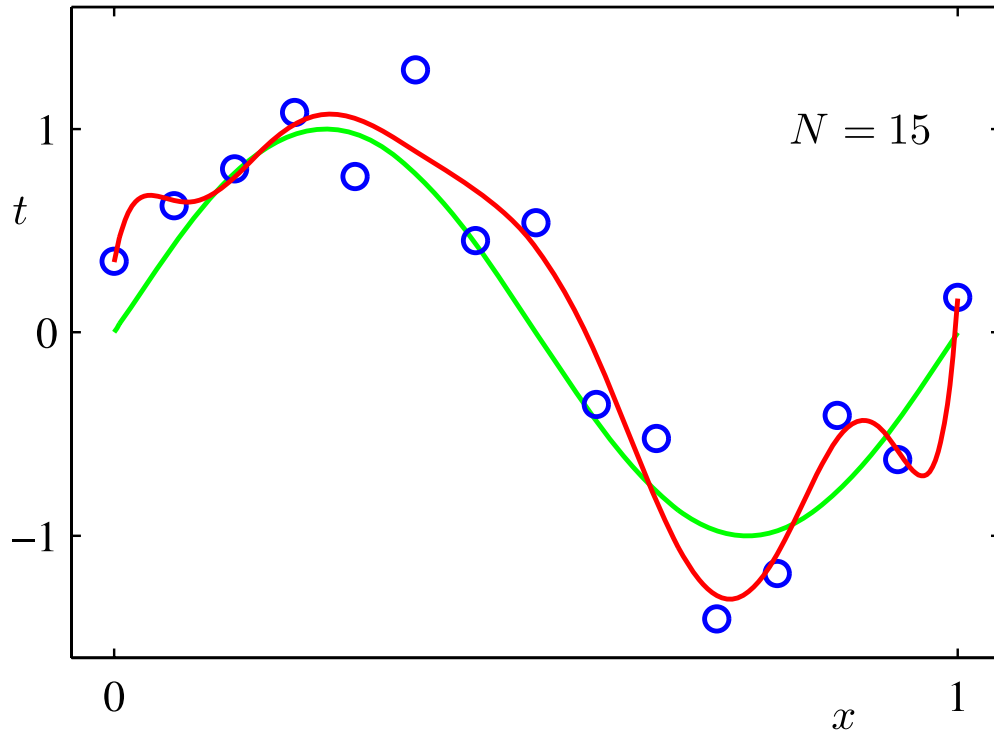
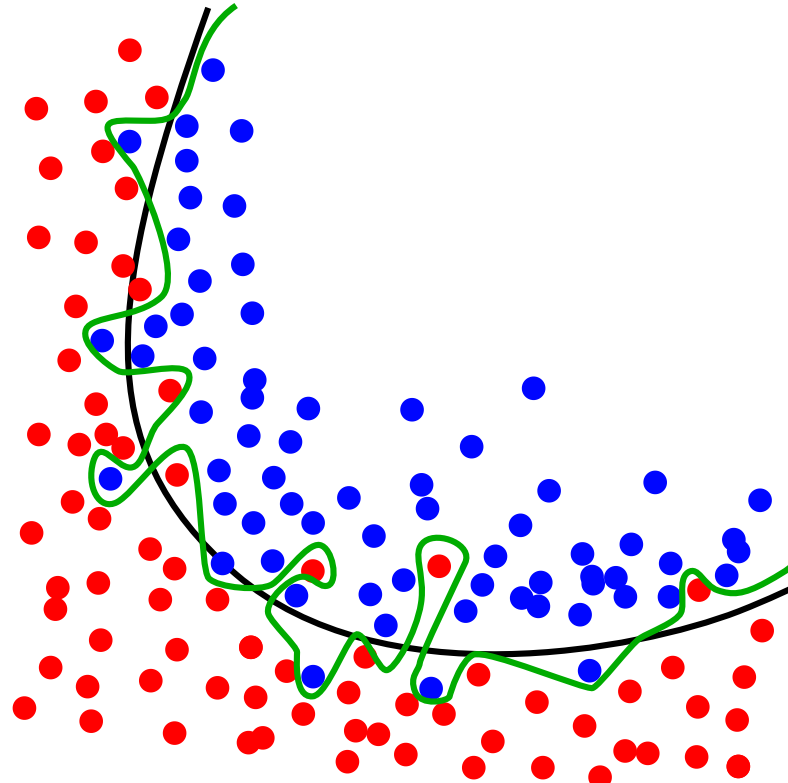


Figure 1.6 of Pattern Recognition and Machine Learning.

Regularization, in a broad sense, is any change that is designed to *reduce generalization error* (but not necessarily its training error) in a machine learning algorithm.

We already saw that **limiting model capacity** can work as regularization.



<https://upload.wikimedia.org/wikipedia/commons/1/19/Overfitting.svg>

L^2 -regularization is one of the oldest regularization techniques, which tries to prefer “simpler” models by endorsing models with **smaller weights**.

Concretely, **L^2 -regularization** (also called **weight decay**) penalizes models with large weights by utilizing the following error function:

$$\frac{1}{2} \sum_{i=1}^N (y(\mathbf{x}_i; \mathbf{w}) - t_i)^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2.$$

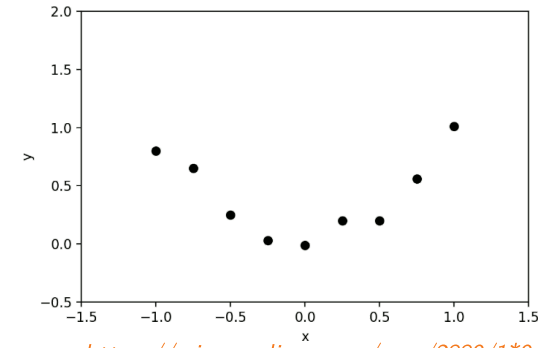
Note that the L^2 -regularization is usually not applied to the *bias*, only to the “proper” weights, because we cannot really overfit via the bias. Also, without penalizing the bias, linear regression with L^2 -regularization is invariant to shifts (i.e., adding a constant to all the targets results in the same solution, only with the bias increased by that constant; if the bias were penalized, this would not be true).

For simplicity, we will not explicitly exclude the bias from the L^2 -regularization penalty in the slides (several textbooks also take the same approach).

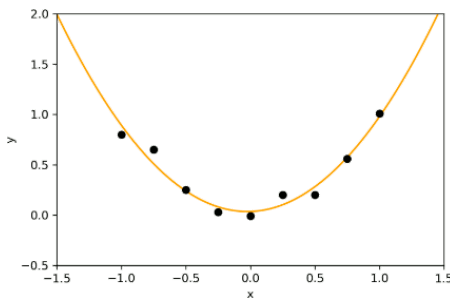
L2 Regularization

One way to look at L^2 -regularization is that it promotes smaller changes of the model (the gradient of linear regression with respect to the inputs are exactly the weights, i.e., $\nabla_x y(\mathbf{x}; \mathbf{w}) = \mathbf{w}$).

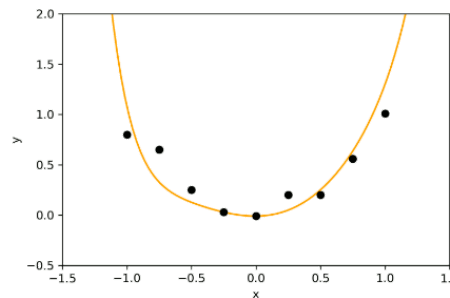
Considering the data points on the right, we present mean squared errors and L^2 norms of the weights for three linear regression models:



https://miro.medium.com/max/2880/1*0-fsK9RkqL3rogo2SnZPCg.png

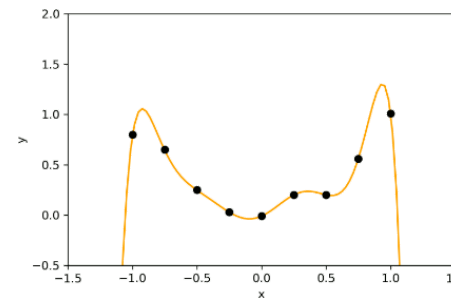


(a) #params = 3
MSE = 0.006
L2 norm = 0.90
L1 norm = 0.98



(b) #params = 9
MSE = 0.035
L2 norm = 1.06
L1 norm = 2.32

https://miro.medium.com/max/2880/1*DVfYChNDMNIS_7CVq2PhSQ.png



(c) #params = 9
MSE = 0
L2 norm = 32.69
L1 norm = 70.03

Figure a: $\hat{y} = 0.04 + 0.04x + 0.9x^2$

Figure b: $\hat{y} = -0.01 + 0.01x + 0.8x^2 + 0.5x^3 - 0.1x^4 - 0.1x^5 + 0.3x^6 - 0.3x^7 + 0.2x^8$

Figure c: $\hat{y} = -0.01 + 0.57x + 2.67x^2 - 4.08x^3 - 12.25x^4 + 7.41x^5 + 24.87x^6 - 3.79x^7 - 14.38x^8$

https://miro.medium.com/max/2880/1*UolRIKXikCz7SFsPFSZrYQ.png

L2 Regularization

The effect of L^2 -regularization can be seen as limiting the *effective capacity* of the model.

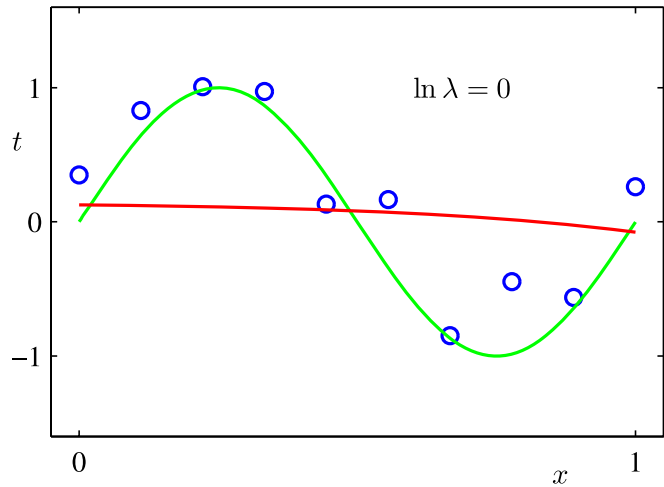
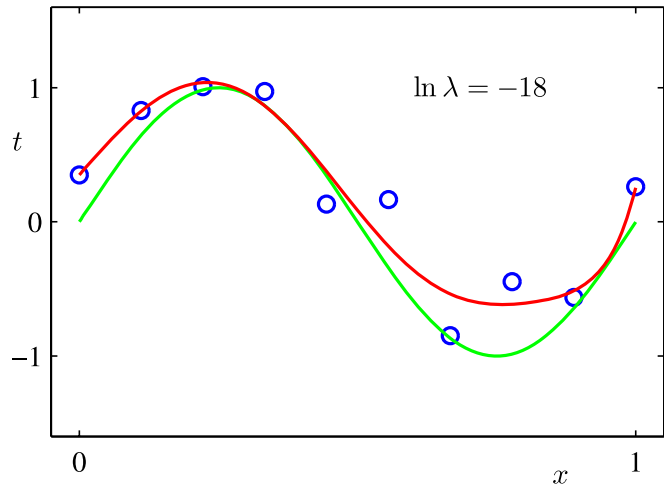


Figure 1.7 of Pattern Recognition and Machine Learning.

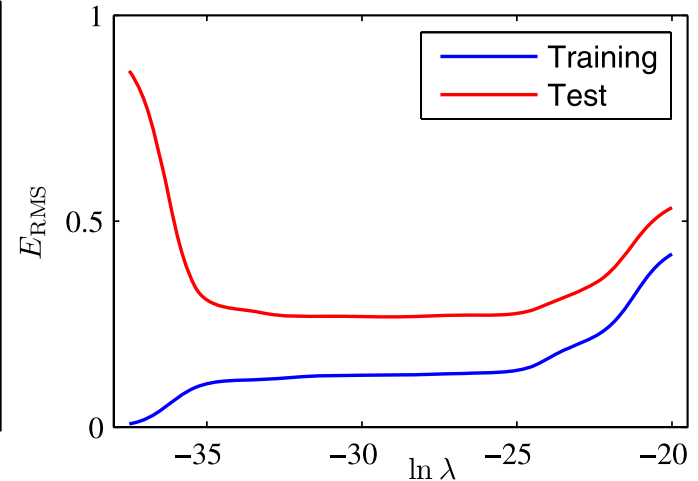


Figure 1.8 of Pattern Recognition and Machine Learning.

Regularizing Linear Regression

In a matrix form, the regularized *sum of squares error* for linear regression amounts to

$$\frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{t}\|^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2.$$

When repeating the same calculation as in the unregularized case, we arrive at

$$(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})\mathbf{w} = \mathbf{X}^T \mathbf{t},$$

where \mathbf{I} is an identity matrix.

Input: Dataset $(\mathbf{X} \in \mathbb{R}^{N \times D}, \mathbf{t} \in \mathbb{R}^N)$, constant $\lambda \in \mathbb{R}^+$.

Output: Weights $\mathbf{w} \in \mathbb{R}^D$ minimizing MSE of regularized linear regression.

- $\mathbf{w} \leftarrow (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{t}.$

Note that the matrix $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$ is always regular for $\lambda > 0$ (you can show that the matrix is positive definite), so another effect of L^2 -regularization is that the inverse always exists.

Choosing Hyperparameters

Hyperparameters are not adapted by the learning algorithm itself.

Usually, a **validation set** or **development set** is used to estimate the generalization error, allowing us to update hyperparameters accordingly. If there is not enough data (well, there is **always** not enough data), more sophisticated approaches can be used.

So far, we have seen two hyperparameters, M and λ .

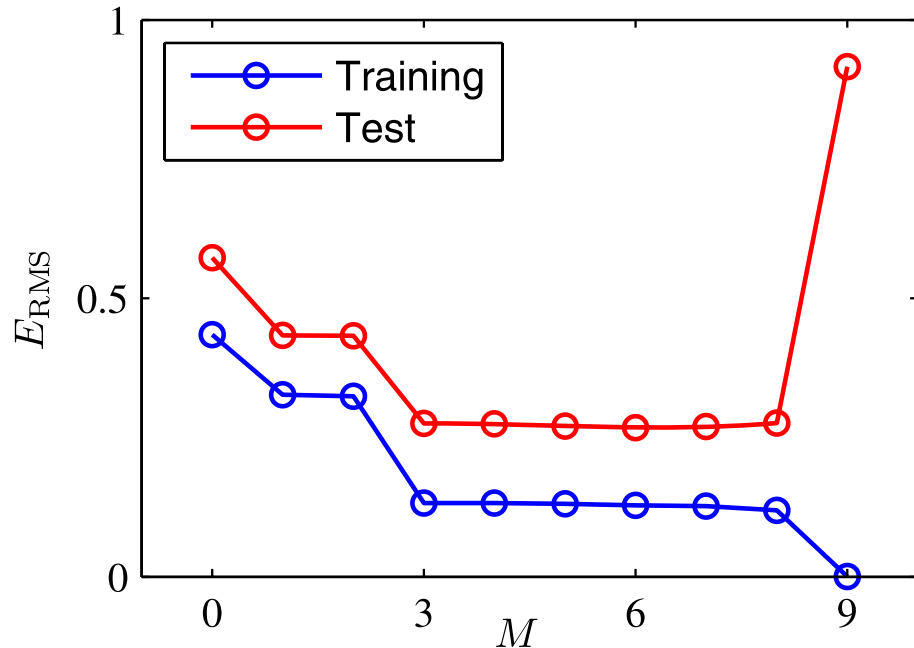


Figure 1.5 of Pattern Recognition and Machine Learning.

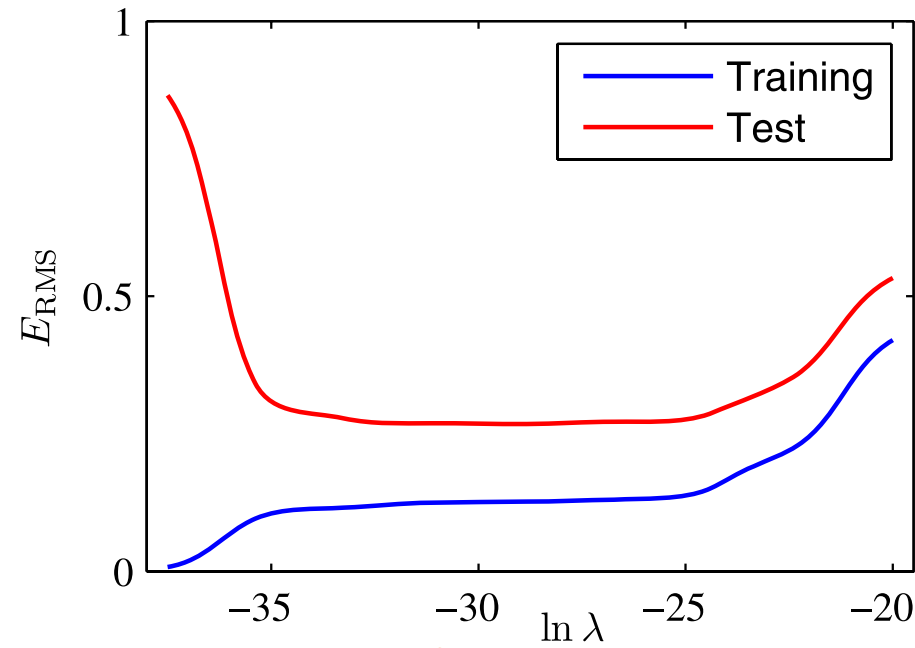


Figure 1.8 of Pattern Recognition and Machine Learning.

When training a linear regression model, we minimized the *sum of squares* error function by computing its gradient (partial derivatives with respect to all weights) and setting it to zero, arriving at the following equation for optimal weights:

$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{t}.$$

If $\mathbf{X}^T \mathbf{X}$ is regular, we can invert it and compute the weights as $\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}$.

It can be proven (see next slide) that $\text{rank}(\mathbf{X}) = \text{rank}(\mathbf{X}^T \mathbf{X})$. Therefore, the matrix $\mathbf{X}^T \mathbf{X} \in \mathbb{R}^{D \times D}$ is regular if and only if \mathbf{X} has rank D , which is equivalent to the columns of \mathbf{X} being linearly independent.

Linear Regression Solution Always Exists

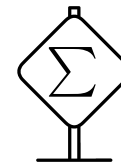
We now show that the solution of $\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{t}$ always exists.

Recall that the rank-nullity theorem states that for a matrix $\mathbf{A} \in \mathbb{R}^{V \times W}$,

$$\text{rank}(\mathbf{A}) + \text{nullity}(\mathbf{A}) \stackrel{\text{def}}{=} \dim(\text{im}(\mathbf{A})) + \dim(\text{ker}(\mathbf{A})) = W.$$

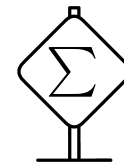
Our goal is to show that $\text{im}(\mathbf{X}^T \mathbf{X}) = \text{im}(\mathbf{X}^T)$. Then the solution would always exist, because for any \mathbf{t} , $\mathbf{X}^T \mathbf{t} \in \text{im}(\mathbf{X}^T \mathbf{X})$.

- We first show that $\text{ker}(\mathbf{X}^T \mathbf{X}) = \text{ker}(\mathbf{X})$.
 - If $\mathbf{X} \mathbf{t} = 0$, then also $\mathbf{X}^T \mathbf{X} \mathbf{t} = 0$, so $\text{ker}(\mathbf{X}^T \mathbf{X}) \supseteq \text{ker}(\mathbf{X})$.
 - If $\mathbf{X}^T \mathbf{X} \mathbf{t} = 0$, then also $\mathbf{t}^T \mathbf{X}^T \mathbf{X} \mathbf{t} = 0$. Therefore $(\mathbf{X} \mathbf{t})^T (\mathbf{X} \mathbf{t}) = 0$, which implies $\mathbf{X} \mathbf{t} = 0$, resulting in $\text{ker}(\mathbf{X}^T \mathbf{X}) \subseteq \text{ker}(\mathbf{X})$.
- Therefore, the rank-nullity theorem implies that $\text{rank}(\mathbf{X}^T \mathbf{X}) = \text{rank}(\mathbf{X}) = \text{rank}(\mathbf{X}^T)$.
- Finally, it is easy to see that $\text{im}(\mathbf{X}^T \mathbf{X}) \subseteq \text{im}(\mathbf{X}^T)$, which together with the rank equality proves the required equation $\text{im}(\mathbf{X}^T \mathbf{X}) = \text{im}(\mathbf{X}^T)$.



SVD Solution of Linear Regression

Now consider the case that $\mathbf{X}^T \mathbf{X}$ is singular. We already know that $\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{t}$ is solvable, but it does not have a unique solution (it has many solutions). Our goal in this case will be to find the \mathbf{w} with the minimum $\|\mathbf{w}\|$ fulfilling the equation.



We now consider *singular value decomposition* (SVD) of \mathbf{X} , writing $\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$, where

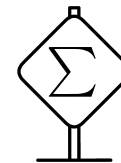
- $\mathbf{U} \in \mathbb{R}^{N \times N}$ is an orthogonal matrix, i.e., $\mathbf{u}_i^T \mathbf{u}_j = [i = j] \Leftrightarrow \mathbf{U}^T \mathbf{U} = \mathbf{I} \Leftrightarrow \mathbf{U}^{-1} = \mathbf{U}^T$,
- $\mathbf{\Sigma} \in \mathbb{R}^{N \times D}$ is a diagonal matrix,
- $\mathbf{V} \in \mathbb{R}^{D \times D}$ is again an orthogonal matrix.

Assuming the diagonal matrix $\mathbf{\Sigma}$ has a rank r , we have

$$\mathbf{\Sigma} = \begin{bmatrix} \mathbf{\Sigma}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix},$$

where $\mathbf{\Sigma}_r \in \mathbb{R}^{r \times r}$ is a regular diagonal matrix. Denoting \mathbf{U}_r and \mathbf{V}_r the matrices of first r columns of \mathbf{U} and \mathbf{V} , respectively, we can write $\mathbf{X} = \mathbf{U}_r \mathbf{\Sigma}_r \mathbf{V}_r^T$.

Using the decomposition $\mathbf{X} = \mathbf{U}_r \boldsymbol{\Sigma}_r \mathbf{V}_r^T$, we can rewrite the goal equation as



$$(\mathbf{V}_r \boldsymbol{\Sigma}_r^T \mathbf{U}_r^T) (\mathbf{U}_r \boldsymbol{\Sigma}_r \mathbf{V}_r^T) \mathbf{w} = (\mathbf{V}_r \boldsymbol{\Sigma}_r^T \mathbf{U}_r^T) \mathbf{t}.$$

The transposition of an orthogonal matrix is its inverse. Therefore, our submatrix \mathbf{U}_r fulfills $\mathbf{U}_r^T \mathbf{U}_r = \mathbf{I}$, because $\mathbf{U}_r^T \mathbf{U}_r$ is the top left submatrix of $\mathbf{U}^T \mathbf{U}$. Analogously, $\mathbf{V}_r^T \mathbf{V}_r = \mathbf{I}$.

We therefore simplify the goal equation to

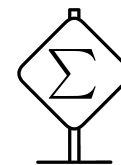
$$\mathbf{V}_r^T \mathbf{V}_r \boldsymbol{\Sigma}_r^T \mathbf{U}_r^T \mathbf{U}_r \boldsymbol{\Sigma}_r \mathbf{V}_r^T \mathbf{w} = \mathbf{V}_r^T \mathbf{V}_r \boldsymbol{\Sigma}_r^T \mathbf{U}_r^T \mathbf{t}.$$

$$\boldsymbol{\Sigma}_r^T \boldsymbol{\Sigma}_r \mathbf{V}_r^T \mathbf{w} = \boldsymbol{\Sigma}_r^T \mathbf{U}_r^T \mathbf{t}$$

Because the diagonal matrix $\boldsymbol{\Sigma}_r = \boldsymbol{\Sigma}_r^T$ is regular, we can divide by it and obtain

$$\mathbf{V}_r^T \mathbf{w} = \boldsymbol{\Sigma}_r^{-1} \mathbf{U}_r^T \mathbf{t}.$$

We have $\mathbf{V}_r^T \mathbf{w} = \boldsymbol{\Sigma}_r^{-1} \mathbf{U}_r^T \mathbf{t}$. If the original matrix $\mathbf{X}^T \mathbf{X}$ was regular, then $r = D$ and \mathbf{V}_r is a square regular orthogonal matrix, in which case $\mathbf{w} = \mathbf{V}_r \boldsymbol{\Sigma}_r^{-1} \mathbf{U}_r^T \mathbf{t}$.



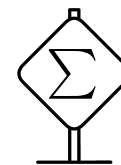
Let $\boldsymbol{\Sigma}^+ \in \mathbb{R}^{D \times N}$ denote the diagonal matrix with

$$\Sigma_{i,i}^+ \stackrel{\text{def}}{=} \begin{cases} \Sigma_{i,i}^{-1} & \text{if } \Sigma_{i,i} \neq 0, \\ 0 & \text{otherwise.} \end{cases}$$

Using this notation, we can rewrite \mathbf{w} for the $r = D$ case to $\mathbf{w} = \mathbf{V} \boldsymbol{\Sigma}^+ \mathbf{U}^T \mathbf{t}$.

Now if $r < D$, $\mathbf{V}_r^T \mathbf{w} = \mathbf{y}$ is undetermined and has infinitely many solutions. To find the one with the smallest norm $\|\mathbf{w}\|$, consider the full product $\mathbf{V}^T \mathbf{w}$. Because \mathbf{V} is orthogonal, $\|\mathbf{V}^T \mathbf{w}\| = \|\mathbf{w}\|$, and it is sufficient to find \mathbf{w} with the smallest $\|\mathbf{V}^T \mathbf{w}\|$. We know that the first r elements of $\mathbf{V}^T \mathbf{w}$ are fixed by the above equation – therefore, the smallest $\|\mathbf{V}^T \mathbf{w}\|$ can be obtained by setting the last $D - r$ elements to zero. Finally, note that $\boldsymbol{\Sigma}^+ \mathbf{U}^T \mathbf{t}$ is exactly $\boldsymbol{\Sigma}_r^{-1} \mathbf{U}_r^T \mathbf{t}$ padded with $D - r$ zeros, which yields the same solution $\mathbf{w} = \mathbf{V} \boldsymbol{\Sigma}^+ \mathbf{U}^T \mathbf{t}$.

The solution to a linear regression with *sum of squares* error function is tightly connected to matrix pseudoinverses. If a matrix \mathbf{X} is singular or rectangular, it does not have an exact inverse, and $\mathbf{X}\mathbf{w} = \mathbf{b}$ does not have an exact solution.



However, we can consider the so-called *Moore-Penrose pseudoinverse*

$$\mathbf{X}^+ \stackrel{\text{def}}{=} \mathbf{V}\mathbf{\Sigma}^+\mathbf{U}^T$$

to be the closest approximation to an inverse, in the sense that we can find the best solution (with the smallest MSE) to the equation $\mathbf{X}\mathbf{w} = \mathbf{b}$ by setting $\mathbf{w} = \mathbf{X}^+\mathbf{b}$.

Alternatively, we can define the pseudoinverse of a matrix \mathbf{X} as

$$\mathbf{X}^+ = \arg \min_{\mathbf{Y} \in \mathbb{R}^{D \times N}} \|\mathbf{X}\mathbf{Y} - \mathbf{I}_N\|_F = \arg \min_{\mathbf{Y} \in \mathbb{R}^{D \times N}} \|\mathbf{Y}\mathbf{X} - \mathbf{I}_D\|_F$$

which can be verified to be the same as our SVD formula.

A random variable \mathbf{x} is a result of a random process, and it can be either discrete or continuous.

Probability Distribution

A probability distribution describes how likely are the individual values that a random variable can take.

The notation $\mathbf{x} \sim P$ stands for a random variable \mathbf{x} having a distribution P .

For discrete variables, the probability that \mathbf{x} takes a value x is denoted as $P(x)$ or explicitly as $P(\mathbf{x} = x)$. All probabilities are nonnegative, and the sum of the probabilities of all possible values of \mathbf{x} is $\sum_x P(\mathbf{x} = x) = 1$.

For continuous variables, the probability that the value of \mathbf{x} lies in the interval $[a, b]$ is given by $\int_a^b p(x) dx$, where $p(x)$ is the *probability density function*, which is always nonnegative and integrates to 1 over the range of all values of \mathbf{x} .

Joint, Conditional, Marginal Probability

For two random variables, a **joint probability distribution** is a distribution of all possible pairs of outputs (and analogously for more than two):

$$P(\mathbf{x} = x_2, y = y_1).$$

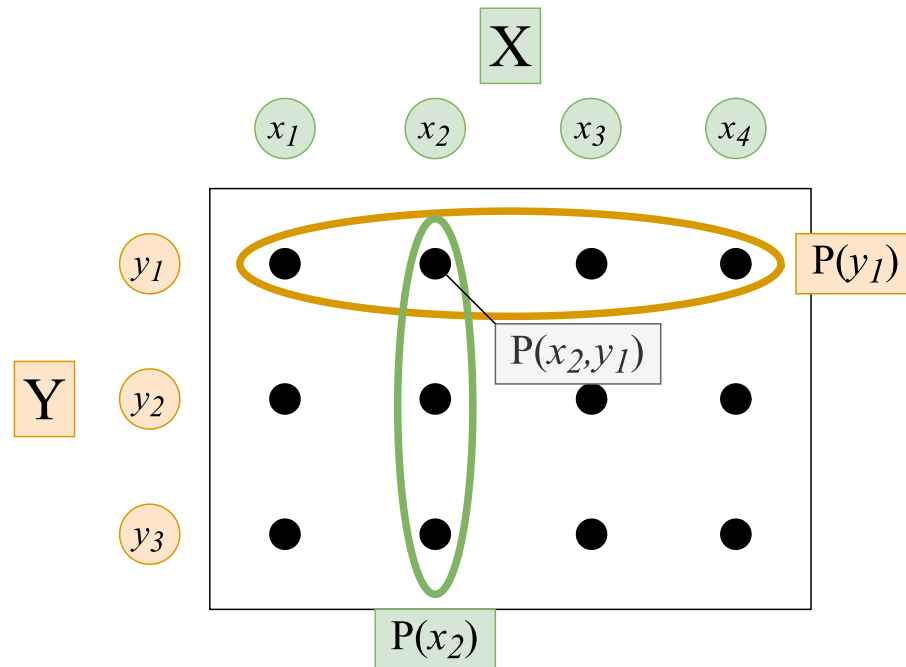
Marginal distribution is a distribution of one (or a subset) of the random variables and can be obtained by summing over the other variable(s):

$$P(\mathbf{x} = x_2) = \sum_y P(\mathbf{x} = x_2, y = y).$$

Conditional distribution is a distribution of one (or a subset) of the random variables, given that another event has already occurred:

$$P(\mathbf{x} = x_2 | y = y_1) = P(\mathbf{x} = x_2, y = y_1) / P(y = y_1).$$

If $P(\mathbf{x}, y) = P(\mathbf{x}) \cdot P(y)$ or $P(\mathbf{x}|y) = P(\mathbf{x})$, random variables \mathbf{x} and y are **independent**.



Expectation

The expectation of a function $f(x)$ with respect to a discrete probability distribution $P(x)$ is defined as:

$$\mathbb{E}_{x \sim P}[f(x)] \stackrel{\text{def}}{=} \sum_x P(x) f(x).$$

For continuous variables, the expectation is computed as:

$$\mathbb{E}_{x \sim p}[f(x)] \stackrel{\text{def}}{=} \int_x p(x) f(x) dx.$$

If the random variable is obvious from context, we can write only $\mathbb{E}_P[x]$, $\mathbb{E}_x[x]$, or even $\mathbb{E}[x]$.

Expectation is linear, i.e., for constants $\alpha, \beta \in \mathbb{R}$:

$$\mathbb{E}_x[\alpha f(x) + \beta g(x)] = \alpha \mathbb{E}_x[f(x)] + \beta \mathbb{E}_x[g(x)].$$

Variance

Variance measures how much the values of a random variable differ from its mean $\mathbb{E}[x]$.

$$\text{Var}(x) \stackrel{\text{def}}{=} \mathbb{E} \left[(x - \mathbb{E}[x])^2 \right], \text{ or more generally,}$$

$$\text{Var}_{x \sim P}(f(x)) \stackrel{\text{def}}{=} \mathbb{E} \left[(f(x) - \mathbb{E}[f(x)])^2 \right].$$

It is easy to see that

$$\text{Var}(x) = \mathbb{E} \left[x^2 - 2x \cdot \mathbb{E}[x] + (\mathbb{E}[x])^2 \right] = \mathbb{E} [x^2] - (\mathbb{E}[x])^2,$$

because $\mathbb{E} [2x \cdot \mathbb{E}[x]] = 2(\mathbb{E}[x])^2$.

Variance is connected to $\mathbb{E}[x^2]$, the **second moment** of a random variable – it is in fact a **centered** second moment.

Estimators and Bias

An **estimator** is a rule for computing an estimate of a given value, often an expectation of some random value(s).

For example, we might estimate *mean* of a random variable by sampling a value according to its probability distribution.

Bias of an estimator is the difference between the expected value of the estimator and the true value being estimated:

$$\textit{estimator bias} \stackrel{\text{def}}{=} \mathbb{E}_{\textit{estimator}}[\textit{estimate}] - \textit{true estimated value}.$$

If the bias is zero, we call the estimator **unbiased**; otherwise, we call it **biased**.

As an example, consider estimating $\mathbb{E}_P[f(x)]$ by generating a single sample x from P and returning $f(x)$. Such an estimate is unbiased, because $\mathbb{E}[\textit{estimate}] = \mathbb{E}_P[f(x)]$, which is exactly the true estimated value.

Estimators and Bias

If we have a sequence of estimates, it might also happen that the bias converges to zero. Consider the well-known sample estimate of variance. Given independent and identically distributed random variables x_1, \dots, x_n , we might estimate the mean and variance as

$$\hat{\mu} = \frac{1}{n} \sum_i x_i, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_i (x_i - \hat{\mu})^2.$$

Such an estimate is biased, because $\mathbb{E}[\hat{\sigma}^2] = (1 - \frac{1}{n})\sigma^2$, but the bias converges to zero with increasing n .

Also, an unbiased estimator does not necessarily have a small variance – in some cases, it can have a large variance, so a biased estimator with a smaller variance might be preferred.

Gradient Descent

Sometimes it is more practical to search for the best model weights in an iterative/incremental/sequential fashion. Either because there is too much data, or the direct optimization is not feasible.

Assuming we are minimizing an error function

$$\arg \min_{\mathbf{w}} E(\mathbf{w}),$$

we may use *gradient descent*:

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \nabla_{\mathbf{w}} E(\mathbf{w})$$

The constant α is called a **learning rate** and specifies the “length” of a step we perform in every iteration of the gradient descent.

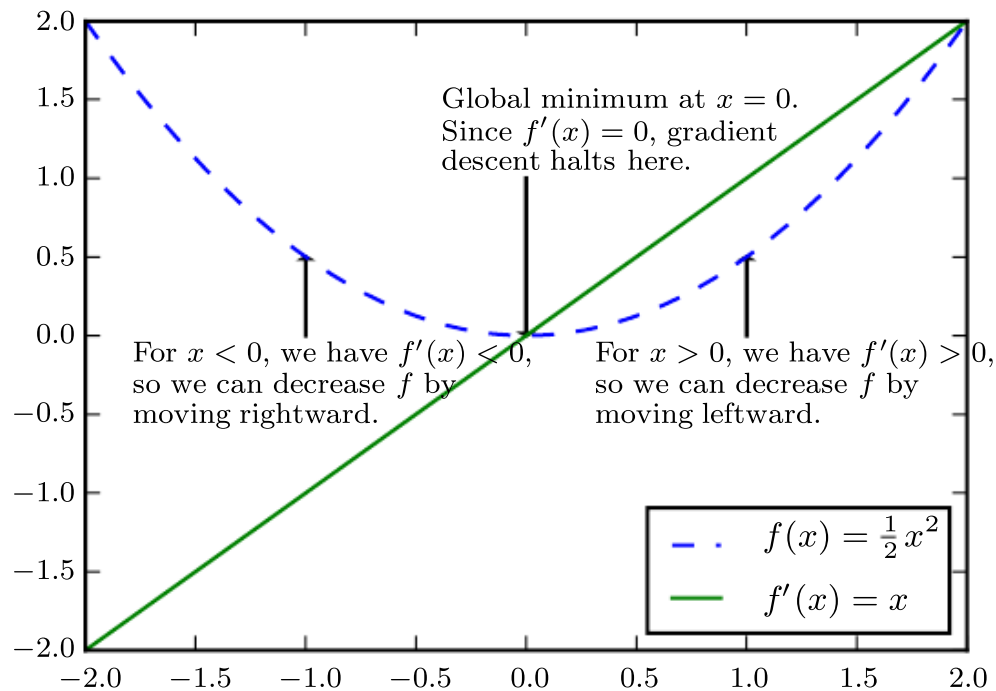


Figure 4.1 of "Deep Learning" book, <https://www.deeplearningbook.org>

Gradient Descent Variants

Let $\mathbf{X} \in \mathbb{R}^{N \times D}$, $\mathbf{t} \in \mathbb{R}^N$ be the training data, and denote $\hat{p}_{\text{data}}(\mathbf{x}, t) \stackrel{\text{def}}{=} \frac{|\{i: (\mathbf{x}, t) = (\mathbf{x}_i, t_i)\}|}{N}$.

Assume that the error function can be computed as an expectation over the dataset:

$$E(\mathbf{w}) = \mathbb{E}_{(\mathbf{x}, t) \sim \hat{p}_{\text{data}}} L(y(\mathbf{x}; \mathbf{w}), t), \quad \text{so that} \quad \nabla_{\mathbf{w}} E(\mathbf{w}) = \mathbb{E}_{(\mathbf{x}, t) \sim \hat{p}_{\text{data}}} \nabla_{\mathbf{w}} L(y(\mathbf{x}; \mathbf{w}), t).$$

- **(Standard/Batch) Gradient Descent:** We use all training data to compute $\nabla_{\mathbf{w}} E(\mathbf{w})$.
- **Stochastic (or Online) Gradient Descent:** We estimate $\nabla_{\mathbf{w}} E(\mathbf{w})$ using a single random example from the training data. Such an estimate is unbiased, but very noisy.

$$\nabla_{\mathbf{w}} E(\mathbf{w}) \approx \nabla_{\mathbf{w}} L(y(\mathbf{x}; \mathbf{w}), t) \quad \text{for a randomly chosen } (\mathbf{x}, t) \text{ from } \hat{p}_{\text{data}}.$$

- **Minibatch SGD:** Trade-off between gradient descent and SGD – the expectation in $\nabla_{\mathbf{w}} E(\mathbf{w})$ is estimated using B random independent examples from the training data.

$$\nabla_{\mathbf{w}} E(\mathbf{w}) \approx \frac{1}{B} \sum_{i=1}^B \nabla_{\mathbf{w}} L(y(\mathbf{x}_i; \mathbf{w}), t_i) \quad \text{for a randomly chosen } (\mathbf{x}_i, t_i) \text{ from } \hat{p}_{\text{data}}.$$

Gradient Descent Convergence

Assume that we perform a stochastic gradient descent, using a sequence of learning rates α_i , and using a noisy estimate $J(\mathbf{w})$ of the real gradient $\nabla_{\mathbf{w}}E(\mathbf{w})$:

$$\mathbf{w}_{i+1} \leftarrow \mathbf{w}_i - \alpha_i J(\mathbf{w}_i).$$

It can be proven (under some reasonable conditions; see Robbins and Monro algorithm, 1951) that if the loss function L is convex and continuous, then SGD converges to the unique optimum almost surely if the sequence of learning rates α_i fulfills the following conditions:

$$\forall i : \alpha_i > 0, \quad \sum_i \alpha_i = \infty, \quad \sum_i \alpha_i^2 < \infty.$$

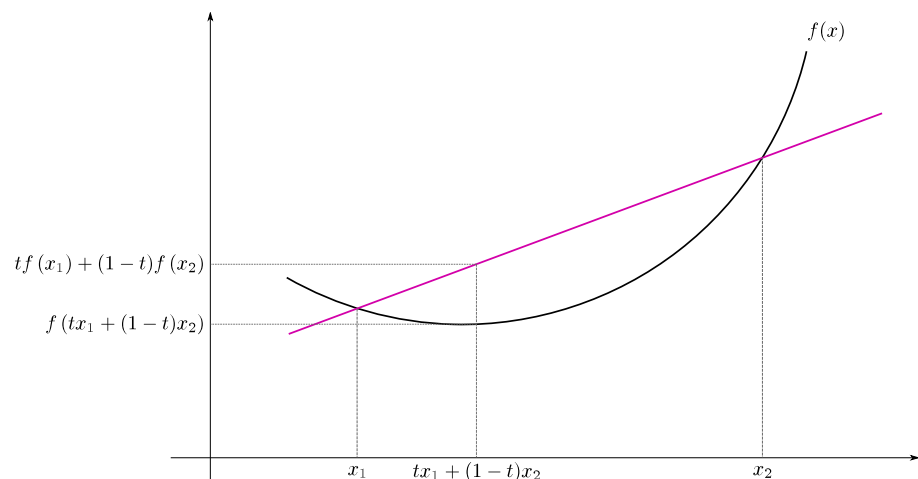
Note that the third condition implies that $\alpha_i \rightarrow 0$.

For nonconvex loss functions, we can get guarantees of converging to a *local* optimum only. However, note that finding the global minimum of an arbitrary function is *at least NP-hard*.

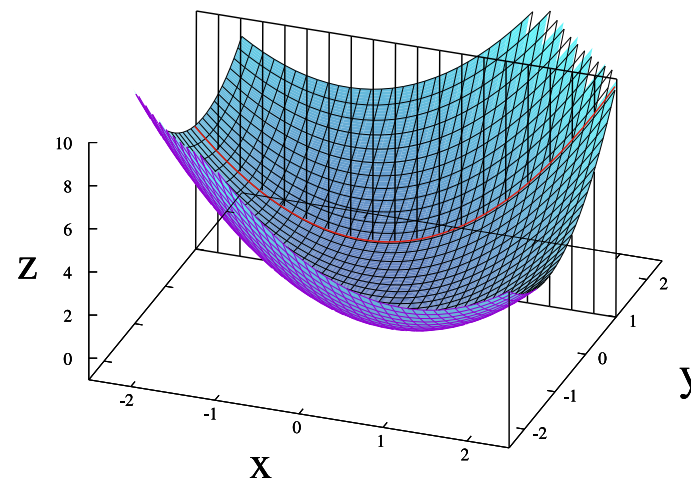
Gradient Descent Convergence

Convex functions mentioned on the previous slide are such that for \mathbf{u}, \mathbf{v} and real $0 \leq t \leq 1$,

$$f(t\mathbf{u} + (1 - t)\mathbf{v}) \leq tf(\mathbf{u}) + (1 - t)f(\mathbf{v}).$$



<https://upload.wikimedia.org/wikipedia/commons/c/c7/ConvexFunction.svg>



https://commons.wikimedia.org/wiki/File:Partial_func_eg.svg

A twice-differentiable function of a single variable is convex iff its second derivative is always nonnegative. (For functions of multiple variables, the Hessian must be positive semi-definite.)

A local minimum of a convex function is always the unique global minimum.

Well-known examples of convex functions are x^2 , e^x , $-\log x$, and also the *sum of squares*.

Solving Linear Regression using SGD

To apply SGD on linear regression, we usually minimize one half of the mean squared error:

$$E(\mathbf{w}) = \mathbb{E}_{(\mathbf{x}, t) \sim \hat{p}_{\text{data}}} \left[\frac{1}{2} (y(\mathbf{x}; \mathbf{w}) - t)^2 \right] = \mathbb{E}_{(\mathbf{x}, t) \sim \hat{p}_{\text{data}}} \left[\frac{1}{2} (\mathbf{x}^T \mathbf{w} - t)^2 \right].$$

If we also include L^2 regularization, we get

$$E(\mathbf{w}) = \mathbb{E}_{(\mathbf{x}, t) \sim \hat{p}_{\text{data}}} \left[\frac{1}{2} (\mathbf{x}^T \mathbf{w} - t)^2 \right] + \frac{\lambda}{2} \|\mathbf{w}\|^2.$$

We then estimate the expectation by a minibatch of examples with indices \mathbb{B} as

$$\frac{1}{|\mathbb{B}|} \sum_{i \in \mathbb{B}} \left(\frac{1}{2} (\mathbf{x}_i^T \mathbf{w} - t_i)^2 \right) + \frac{\lambda}{2} \|\mathbf{w}\|^2,$$

which gives us an estimate of a gradient

$$\nabla_{\mathbf{w}} E(\mathbf{w}) \approx \frac{1}{|\mathbb{B}|} \sum_{i \in \mathbb{B}} \left((\mathbf{x}_i^T \mathbf{w} - t_i) \mathbf{x}_i \right) + \lambda \mathbf{w}.$$

Solving Linear Regression using SGD

The computed gradient allows us to formulate the following algorithm for solving linear regression with minibatch SGD.

Input: Dataset $(\mathbf{X} \in \mathbb{R}^{N \times D}, \mathbf{t} \in \mathbb{R}^N)$, learning rate $\alpha \in \mathbb{R}^+$, L^2 strength $\lambda \in \mathbb{R}$.

Output: Weights $\mathbf{w} \in \mathbb{R}^D$ hopefully minimizing the regularized MSE of a linear regression model.

- $\mathbf{w} \leftarrow \mathbf{0}$ or we initialize \mathbf{w} randomly
- repeat until convergence (or until our patience runs out):
 - sample a minibatch of examples with indices \mathbb{B}
 - either uniformly randomly,
 - or we may want to process all training instances before repeating them, which can be implemented by generating a random permutation and then splitting it into minibatch-sized chunks
 - the most common option; one pass through the data is called an **epoch**

- $$\mathbf{w} \leftarrow \mathbf{w} - \alpha \frac{1}{|\mathbb{B}|} \sum_{i \in \mathbb{B}} ((\mathbf{x}_i^T \mathbf{w} - t_i) \mathbf{x}_i) - \alpha \lambda \mathbf{w}$$

Features

Recall that the *input* instance values are usually the raw observations and are given. However, we might extend them suitably before running a machine learning algorithm, especially if the algorithm is linear or otherwise limited and cannot represent an arbitrary function. Such instance representations are called *features*.

We already saw this in the example from the previous lecture, where even if our training examples were x and t , we performed the linear regression using features (x^0, x^1, \dots, x^M) :

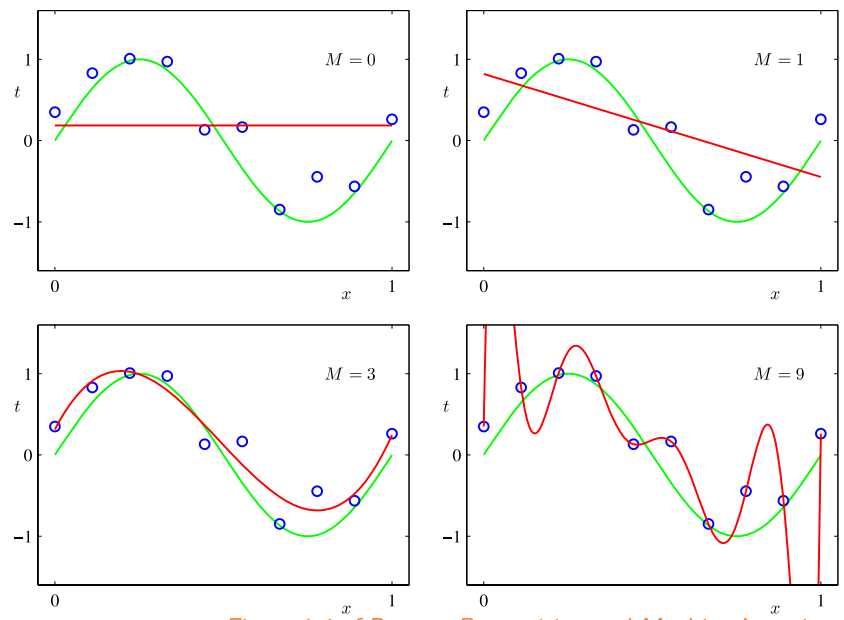


Figure 1.4 of Pattern Recognition and Machine Learning.

Generally, it would be best if the machine learning algorithms would process only the raw inputs. However, many algorithms are capable of representing only a limited set of functions (for example linear ones), and in that case, **feature engineering** plays a major part in the final model performance. Feature engineering is a process of constructing features from raw inputs.

Commonly used features are:

- **polynomial features** of degree p : Given features (x_1, x_2, \dots, x_D) , we might consider *all* products of p input values. Therefore, polynomial features of degree 2 would consist of $x_i^2 \forall i$ and of $x_i x_j \forall i \neq j$.
- **categorical one-hot features**: Assume, for example, that a day in a week is represented in the input as an integer value of 1 to 7, or a breed of a dog is expressed as an integer value of 0 to 366. Using these integral values as an input to linear regression makes little sense – instead, it might be better to learn weights for individual days in a week or for individual dog breeds. We might therefore represent input classes by binary indicators for every class, giving rise to a **one-hot** representation, where an input integral value $0 \leq v < L$ is represented as L binary values, which are all zero except for the v^{th} one, which is one.