



Faculty of Science



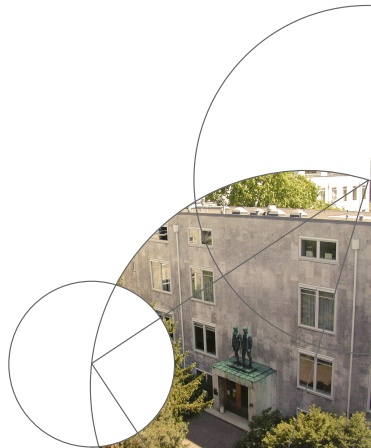
Linear Regression

Machine Learning

Christian Igel
Department of Computer Science

🐦 @christian.igel

Slide 1/42



Outline

- 1 Warm-up: Derivatives and Gradients
- 2 Linear Regression Method
- 3 Deriving Linear Regression
- 4 Non-linear Transformations
- 5 Overfitting
- 6 Regularization
- 7 Summary



Outline

- 1 Warm-up: Derivatives and Gradients
- 2 Linear Regression Method
- 3 Deriving Linear Regression
- 4 Non-linear Transformations
- 5 Overfitting
- 6 Regularization
- 7 Summary



Derivative

Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a real-valued function. If the limit

$$f'(x) = \lim_{\epsilon \rightarrow 0} \frac{f(x + \epsilon) - f(x)}{\epsilon}$$

exists, the function f is differentiable and f' is its derivative.



Differentiation rules

Let f and g be differentiable functions.

Sum. If $h(x) = f(x) + g(x)$ then h is differentiable and $h'(x) = f'(x) + g'(x)$.

Chain. If $h(x) = f(g(x))$ then h is differentiable and $h'(x) = f'(g(x))g'(x)$.

Product. If $h(x) = f(x)g(x)$ then h is differentiable and $h'(x) = f'(x)g(x) + f(x)g'(x)$.

Power. $f(x) = x^r$ with $r \in \mathbb{R}$ has derivative $f'(x) = rx^{r-1}$.

Exponential. For $f(x) = e^x$ we have $f'(x) = e^x$.

Logarithm. For $f(x) = \ln x$ we have $f'(x) = 1/x$.



Important special cases

Constant. Product rule implies $h'(x) = af'(x)$ for $h(x) = af(x)$ with $a \in \mathbb{R}$

Quotient. Product and power give for

$$h(x) = \frac{f(x)}{g(x)}$$

with non-vanishing g (i.e., g is never zero)

$$h'(x) = -\frac{f'(x)g(x) - f(x)g'(x)}{[g(x)]^2}.$$

Reciprocal. Special case, we have $g'(x) = -1/x^2$ for $g(x) = 1/x$ and, more general, if $h(x) = \frac{1}{f(x)}$ for f non-vanishing, then

$$h'(x) = -\frac{f'(x)}{[f(x)]^2}.$$



Partial derivative

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ depend on d variables $(x_1, \dots, x_d)^\top \in \mathbb{R}^d$.

The partial derivative of f with respect to x_i at a point $\mathbf{x} \in \mathbb{R}^d$ is defined as

$$\frac{\partial}{\partial x_i} f(\mathbf{x}) = \lim_{\epsilon \rightarrow 0} \frac{f(\mathbf{x} + \epsilon \mathbf{e}_i) - f(\mathbf{x})}{\epsilon}$$

with \mathbf{e}_i being the i th unit vector (i.e., a vector where the i th component is one and all other components are zero) given that the limit exists. For $d = 1$ we have

$$\frac{\partial}{\partial x} f(x) = f'(x) \ .$$



Gradient

- *Rate of change (directional derivative)* of $f : \mathbb{R}^d \rightarrow \mathbb{R}$ at a point $\mathbf{x} \in \mathbb{R}^d$ when moving in the direction $\mathbf{u} \in \mathbb{R}^d$, $\|\mathbf{u}\| = 1$, is defined as:

$$\nabla_{\mathbf{u}} f(\mathbf{x}) = \lim_{\epsilon \rightarrow 0} \frac{f(\mathbf{x} + \epsilon \mathbf{u}) - f(\mathbf{x})}{\epsilon}$$

- The *gradient*

$$\nabla f(\mathbf{x}) = \left(\frac{\partial f(\mathbf{x})}{\partial x_1}, \frac{\partial f(\mathbf{x})}{\partial x_2}, \dots, \frac{\partial f(\mathbf{x})}{\partial x_d} \right)^T$$

points in the direction $\nabla f(\mathbf{x}) / \|\nabla f(\mathbf{x})\|$ giving maximum rate of change $\|\nabla f(\mathbf{x})\|$.



Chain rule

The *chain rule* for computing the derivative of a composition of two functions,

$$\frac{\partial f(g(x))}{\partial x} = f'(g(x))g'(x)$$

with $f'(x) = \frac{\partial f(x)}{\partial x}$ and $g'(x) = \frac{\partial g(x)}{\partial x}$, can be extended to:

$$\frac{\partial f(g_1(x), g_2(x), \dots, g_d(x))}{\partial x} = \sum_{i=1}^d \frac{\partial f(g_1(x), \dots, g_d(x))}{\partial g_i(x)} \frac{\partial g_i(x)}{\partial x}$$



Outline

- 1 Warm-up: Derivatives and Gradients
- 2 Linear Regression Method
- 3 Deriving Linear Regression
- 4 Non-linear Transformations
- 5 Overfitting
- 6 Regularization
- 7 Summary



The linear regression problem

- The goal is to predict a target variable $\mathbf{y} \in \mathbb{R}^m$ based on given data $\mathbf{x} \in \mathcal{X} = \mathbb{R}^d$. For simplicity, let $m = 1$.
- We build affine linear predictive models

$$y = f(\mathbf{x}) = \sum_{i=1}^d w_i x_i + b = \mathbf{w}^\top \mathbf{x} + b$$

with $\mathbf{w} \in \mathbb{R}^d$ and $b \in \mathbb{R}$.

- For convenience, we define $\tilde{\mathbf{x}}^\top = (x_1, \dots, x_d, 1)$ and $\tilde{\mathbf{w}}^\top = (w_1, \dots, w_d, b)$ and consider the equivalent formulation

$$y = f(\tilde{\mathbf{x}}) = \sum_{i=1}^{d+1} \tilde{w}_i \tilde{x}_i = \tilde{\mathbf{w}}^\top \tilde{\mathbf{x}} .$$



Data matrix

Training data $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ is gathered in data matrix (x_{ij} is the j th component of the i th training pattern, i.e., $[\mathbf{x}_i]_j$)

$$\tilde{\mathbf{X}} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1d} & 1 \\ x_{21} & x_{22} & \dots & x_{2d} & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nd} & 1 \end{pmatrix}$$

Let's omit the tilde to keep notation uncluttered.



Formalizing the goal

- We want to minimize the empirical risk under the squared loss $\ell(y, \hat{y}) = (y - \hat{y})^2$:

$$\frac{1}{n} \sum_{i=1}^n \ell(y_i, \mathbf{w}^\top \mathbf{x}_i) = \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2$$

This quantity is also known as mean-squared-error (MSE).

- Optimal parameters are given by

$$\mathbf{w}^\star = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} ,$$

where $\mathbf{y} = (y_1, \dots, y_n)^\top$.



Linear regression algorithm

Algorithm 1: linear regression

Input: $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\} \in (\mathbb{R}^d \times \mathbb{R})^n$

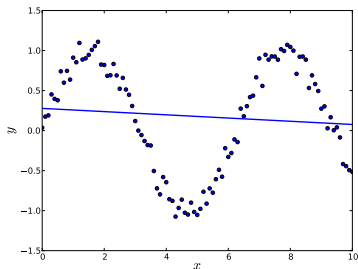
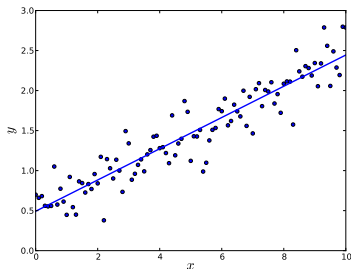
Output: affine linear model

- 1 $\mathbf{y} = (y_1, \dots, y_n)^\top$
- 2 $\tilde{\mathbf{X}} = \begin{pmatrix} x_{11} & \dots & x_{1d} & 1 \\ \vdots & \ddots & \vdots & \vdots \\ x_{n1} & \dots & x_{nd} & 1 \end{pmatrix}$
- 3 $(w_1, \dots, w_d, b)^\top = \left(\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^\top \mathbf{y}$

Result: (\mathbf{w}, b) , defining the model $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$



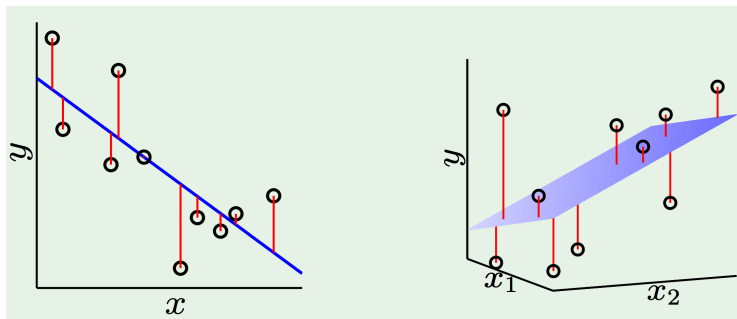
Examples



The right plot is an example of “underfitting”.



Visualizing the error



Abu-Mostafa, Magdon-Ismail, and Lin. *Learning from Data*.
AMLbook, 2012



Outline

- 1 Warm-up: Derivatives and Gradients
- 2 Linear Regression Method
- 3 Deriving Linear Regression**
- 4 Non-linear Transformations
- 5 Overfitting
- 6 Regularization
- 7 Summary



Solution using pseudo-inverse I

Setting the derivative

$$\frac{\partial}{\partial \mathbf{w}} \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 = - \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i) \mathbf{x}_i^\top$$

of $(n/2)$ times the empirical risk to $\mathbf{0}$ yields

$$\mathbf{0}^\top = \sum_{i=1}^n y_i \mathbf{x}_i^\top - \mathbf{w}^\top \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top ,$$

which implies with $\mathbf{y} = (y_1, \dots, y_n)^\top$

$$\mathbf{0}^\top = \sum_{i=1}^n y_i \mathbf{x}_i^\top - \mathbf{w}^\top \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top = \mathbf{y}^\top \mathbf{X} - \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} .$$



Solution using pseudo-inverse II

$$\mathbf{0}^\top = \mathbf{y}^\top \mathbf{X} - \mathbf{w}^\top \mathbf{X}^\top \mathbf{X}$$

implies

$$\mathbf{w}^\top \mathbf{X}^\top \mathbf{X} = \mathbf{y}^\top \mathbf{X}$$

and thus

$$(\mathbf{X}^\top \mathbf{X})^\top \mathbf{w} = \mathbf{X}^\top \mathbf{y} .$$

Assuming $\mathbf{X}^\top \mathbf{X}$ has full rank, we get the maximum likelihood estimate:

$$\mathbf{w}^\star = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

$\mathbf{X}^\dagger = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ is called *Moore-Penrose pseudo-inverse*.

In practice, one does not compute $(\mathbf{X}^\top \mathbf{X})^{-1}$, but considers the *QR*-decomposition of \mathbf{X} and solves a system of linear equations.



Why sum-of-squares?

- Sum-of-squares error has the nice properties of
 - being differentiable and
 - penalizing large deviations from a target value more than proportionally stronger than small deviations.
- There are theoretical reasons . . .



Likelihood function

- Likelihood (function) of the parameters \boldsymbol{w} given training data S is the probability of observing S when the data is generated by h with parameters \boldsymbol{w} .
- Likelihood for i.i.d. S :

$$\prod_{i=1}^N P(Y = y_i \mid X = \boldsymbol{x}_i; h) \text{ or short } \prod_{i=1}^N P(y_i \mid \boldsymbol{x}_i)$$

- Negative log(arithmic) likelihood:

$$-\log \prod_{i=1}^N P(y_i \mid \boldsymbol{x}_i) = -\sum_{i=1}^N \log P(y_i \mid \boldsymbol{x}_i)$$

- Learning principle: Negative log-likelihood minimization



Sum-of-squares and Gaussian noise I

Consider a deterministic real-valued target functions perturbed by zero-mean additive Gaussian noise with variance σ^2 , then we have

$$p(y \mid \mathbf{x}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-f(\mathbf{x}))^2}{2\sigma^2}}.$$

Given $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ the likelihood of the model $h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$ is

$$p((y_1, \dots, y_n)^\top \mid (\mathbf{x}_1, \dots, \mathbf{x}_n), \mathbf{w}, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \mathbf{w}^\top \mathbf{x}_i)^2}{2\sigma^2}}.$$



Sum-of-squares and Gaussian noise II

Let the logarithm turn products into sums:

$$\begin{aligned}\ln p(\mathbf{y} | \mathbf{w}) &= \ln \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \mathbf{w}^\top \mathbf{x}_i)^2}{2\sigma^2}} \\&= \sum_{i=1}^n \ln \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \mathbf{w}^\top \mathbf{x}_i)^2}{2\sigma^2}} \\&= \underbrace{-\frac{n}{2} \ln \sigma^2 - \frac{n}{2} \ln(2\pi)}_{\text{independent of } \mathbf{w}} - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2\end{aligned}$$

Thus, maximizing the likelihood corresponds to minimizing MSE under the assumption of Gaussian noise.

Why Gaussian noise? Central limit theorem, maximum entropy principle



Outline

- 1 Warm-up: Derivatives and Gradients
- 2 Linear Regression Method
- 3 Deriving Linear Regression
- 4 Non-linear Transformations**
- 5 Overfitting
- 6 Regularization
- 7 Summary



Non-linear transformations

Consider a non-linear transformation

$$\Phi : \mathcal{X} \rightarrow \mathcal{Z} .$$

Then we can consider hypotheses

$$h(\mathbf{x}) = \tilde{h}(\Phi(\mathbf{x})) = \tilde{\mathbf{w}}^\top \Phi(\mathbf{x})$$

and do linear regression with the data matrix:

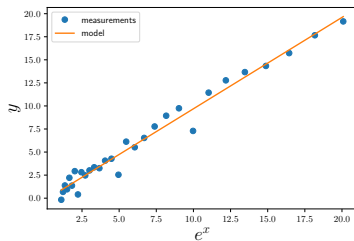
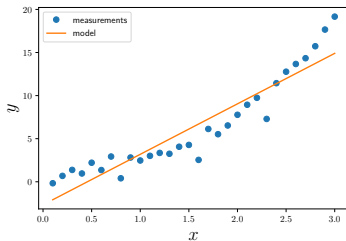
$$\begin{pmatrix} [\Phi(\mathbf{x}_1)]_1 & [\Phi(\mathbf{x}_1)]_2 & \dots & 1 \\ [\Phi(\mathbf{x}_2)]_1 & [\Phi(\mathbf{x}_2)]_2 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ [\Phi(\mathbf{x}_n)]_1 & [\Phi(\mathbf{x}_2)]_2 & \dots & 1 \end{pmatrix}$$



Example I

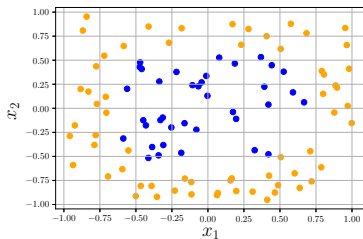
$$\Phi : \mathbb{R} \rightarrow \mathbb{R}$$

$$x \mapsto \exp x$$

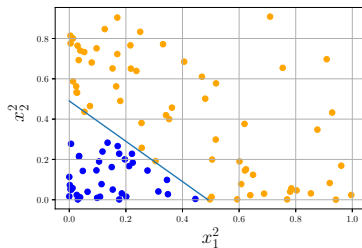


Example II

$$\Phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$$
$$(x_1, x_2) \mapsto (x_1^2, x_2^2)$$



$$(x_1, x_2)^T$$



$$(x_1^2, x_2^2)^T$$



Outline

- 1 Warm-up: Derivatives and Gradients
- 2 Linear Regression Method
- 3 Deriving Linear Regression
- 4 Non-linear Transformations
- 5 Overfitting**
- 6 Regularization
- 7 Summary



Small sample problem

Consider the case $n \leq d$. Then linear regression always finds a model with zero empirical risk by solving the system of linear equations:

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

Even if the y_i are independent of the x_i !

Overfitting

Overfitting occurs if the model faithfully reflects idiosyncrasies of the training data rather than the underlying process that has generated the data. This may happen if the model is too complex/flexible in relation to the amount of available information.

Formal concepts of model complexity and flexibility exist.



Outline

- 1 Warm-up: Derivatives and Gradients
- 2 Linear Regression Method
- 3 Deriving Linear Regression
- 4 Non-linear Transformations
- 5 Overfitting
- 6 Regularization**
- 7 Summary



Empirical risk minimization

So far, we performed *empirical risk minimization*:

$$\operatorname{argmin}_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(y_i, h(\mathbf{x}_i))$$

In the linear regression case, $\ell(y, y') = (y - y')^2$ and $\mathcal{H} = \{h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} \mid \mathbf{w} \in \mathbb{R}^d\}$.

Adding a transformation $\Phi : \mathcal{X} \rightarrow \mathcal{Z}$ gives the hypotheses space

$$\mathcal{H} = \left\{ h(\mathbf{x}) = \mathbf{w}^\top \Phi(\mathbf{x}) \mid \mathbf{w} \in \mathbb{R}^{\tilde{d}} \right\}$$

where \tilde{d} is the feature space dimensionality.



Regularization & model selection

Let $\mathcal{X} = \mathbb{R}$ and

$$\Phi_k : x \mapsto (1, x, x^2, \dots, x^k)^\top .$$

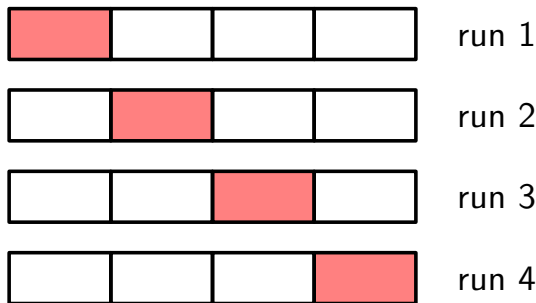
When k is too small, we may underfit. If k and thus \tilde{d} is too large, we may overfit.

How to select the right k ? This is often referred to a *model selection* problem.

→ One way is using (cross-)validation.



Excursus: Cross-validation



C. M. Bishop. *Pattern Recognition and Machine Learning*, Springer, 2006

- K -fold cross-validation (K -CV): Split training data into K (almost) even parts
- Train K times on $K - 1$ parts and test on the held-out part
- Take hyperparameters maximizing mean test performance for building the model using all data



Always remember!

Perhaps most important rule

As a general rule, you must not use the “test data” (i.e., data you use for estimating the final performance of a model) in the model building process at all (neither for training, data normalization, nor hyperparameter selection), because otherwise you may get a biased estimate of the generalization performance of the model.



Two-norm regularization

Let us restrict the hypothesis space to

$$\mathcal{H}_C = \left\{ h(\mathbf{x}) = \mathbf{w}^\top \Phi(\mathbf{x}) \mid \mathbf{w} \in \mathbb{R}^{\tilde{d}} \wedge \|\mathbf{w}\| \leq C \right\}$$

with $C > 0$.

For simplicity, assume $\mathcal{X} = \mathbb{R}$, $\tilde{d} = 1$, and Φ to be Lipschitz with constant α (i.e., $|\Phi(x_1) - \Phi(x_2)| \leq \alpha|x_1 - x_2|$ for $x_1, x_2 \in \mathbb{R}$).

Then the functions in \mathcal{H}_C are Lipschitz with constant $C \cdot \alpha$.

That is, C limits how fast the functions can change – it constrains a notion of complexity of the hypotheses class.



Equivalence

There exists a λ dependent on C such that

$$\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{w}^\top \Phi(\mathbf{x}))^2 \quad \text{s.t.} \quad \|\mathbf{w}\| \leq C$$

is equivalent to:

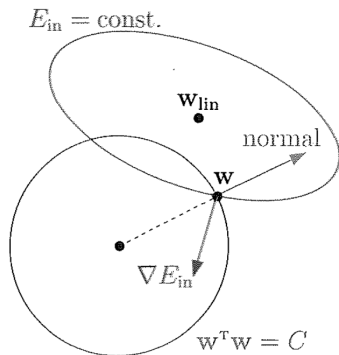
$$\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \left[\frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{w}^\top \Phi(\mathbf{x}))^2 + \lambda \mathbf{w}^\top \mathbf{w} \right]$$

The *regularization parameter* λ can be chosen by (cross-)validation.

Linear regression with two-norm regularization is often called *ridge regression*.



Equivalence proof



E_{in} is the empirical risk, w_{lin} is w^* .

$$\frac{\partial}{\partial w} \left[\frac{1}{n} \sum_{i=1}^n (y_i - f(x))^2 + \lambda w^T w \right] = 0$$

$$\Rightarrow \frac{\partial}{\partial w} \frac{1}{n} \sum_{i=1}^n (y_i - f(x))^2 = -\lambda w^T$$

$$\Rightarrow \nabla E_{\text{in}} = -\lambda w^*$$

Abu-Mostafa, Magdon-Ismail, and Lin. *Learning from Data*. AMLbook, 2012



Solving ridge regression

We had

$$\frac{\partial}{\partial \mathbf{w}} \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 = \mathbf{y}^\top \mathbf{X} - \mathbf{w}^\top \mathbf{X}^\top \mathbf{X}$$

and thus

$$\frac{\partial}{\partial \mathbf{w}} \left[\frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 + \lambda \mathbf{w}^\top \mathbf{w} \right] = \mathbf{y}^\top \mathbf{X} - \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{w}^\top$$

which is solved by:

$$\mathbf{w}^\star = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

What happens to \mathbf{w}^\star for $\lambda \rightarrow \infty$?



One-norm regularization

The one-norm can be used instead of the two-norm:

$$\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^{\tilde{d}}} \left[\frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{w}^T \Phi(\mathbf{x}))^2 + \underbrace{\lambda |\mathbf{w}|}_{\sum_{j=1}^{\tilde{d}} |w_j|} \right]$$

Linear regression with one-norm regularization is referred to as *LASSO* (least absolute shrinkage and selection operator).



Sparsity

One-norm regularization leads to sparser models (i.e., models with several coefficients $w_i = 0$).

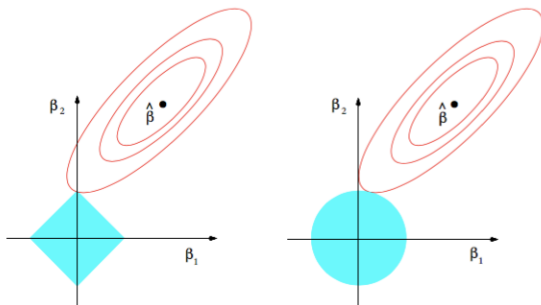


FIGURE 3.11. Estimation picture for the lasso (left) and ridge regression (right). Shown are contours of the error and constraint functions. The solid blue areas are the constraint regions $|\beta_1| + |\beta_2| \leq t$ and $\beta_1^2 + \beta_2^2 \leq t^2$, respectively, while the red ellipses are the contours of the least squares error function.

Outline

- 1 Warm-up: Derivatives and Gradients
- 2 Linear Regression Method
- 3 Deriving Linear Regression
- 4 Non-linear Transformations
- 5 Overfitting
- 6 Regularization
- 7 Summary



Summary

- Regression means building predictive models for real-valued target patterns.
- Linear regression is a baseline method for regression tasks.
- Too complex/flexible models can lead to overfitting.
- However, the affine linear models built by standard linear regression are not very flexible and rather tend to “underfit” the data. Therefore non-linear methods are needed.
- Non-linear transformations combined with linear regression lead to non-linear models.
- Regularization can prevent overfitting.
- 1-norm regularization leads to sparsity.

