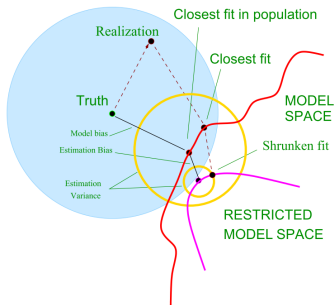


# Einführung in das Statistische Lernen

## Introduction to Regularization



### Learning goals

- Understand why overfitting happens
- Know how overfitting can be avoided
- Know regularized empirical risk minimization

# Motivation for Regularization

# EXAMPLE: OVERFITTING

- Assume we want to predict the daily maximum **ozone level** in LA given a data set containing 50 observations.
- The data set contains 12 features describing time conditions (e.g., weekday, month), the weather (e.g., temperature at different weather stations, humidity, wind speed) or geographic variables (e.g., the pressure gradient).
- We fit a linear regression model using **all** of the features

$$f(\mathbf{x} \mid \boldsymbol{\theta}) = \boldsymbol{\theta}^T \mathbf{x} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_{12} x_{12}$$

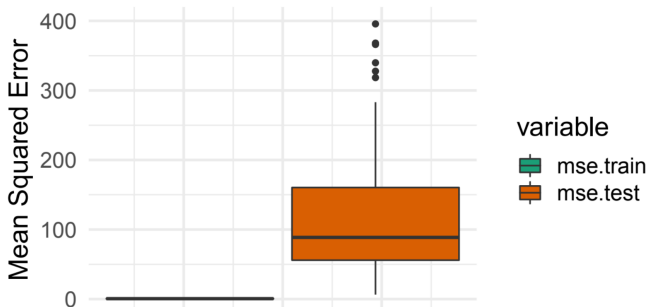
with the  $L2$  loss.

- We evaluate the performance with 10 times 10-fold CV.

We use (a subset of) the `Ozone` data set from the `mlbench` package. This way, we artificially create a “high-dimensional” dataset by reducing the number of observations drastically while keeping the number of features fixed.

# EXAMPLE: OVERFITTING

While our model fits the training data almost perfectly (left), it generalizes poorly to new test data (right). We overfitted.



# AVOID OVERFITTING

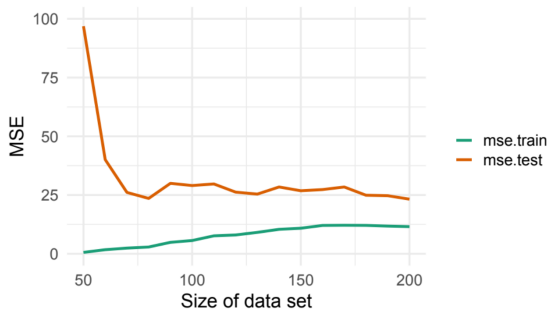
Why can **overfitting** happen? And how to avoid it?

- ❶ Not enough data  
→ collect **more data**
- ❷ Data is noisy  
→ collect **better data** (reduce noise)
- ❸ Models are too complex  
→ use **less complex models**
- ❹ Aggressive loss optimization  
→ **optimize less**

# AVOID OVERFITTING

## Approach 1: Collect more data

We explore our results for increased dataset size by 10 times 10-fold CV. The fit worsens slightly, but the test error decreases.



Good idea, but often not feasible in practice.

# AVOID OVERFITTING

## Approach 3: Reduce complexity

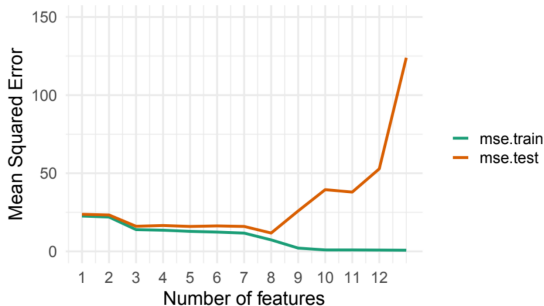
We try the simplest model we can think of: the constant model. For the  $L2$  loss, the optimal constant model is

$$f(\mathbf{x} \mid \theta) = \frac{1}{n} \sum_{i=1}^n y^{(i)}$$

We then increase the complexity of the model step-by-step by adding one feature at a time.

# AVOID OVERFITTING

We can control the complexity of the model by including/excluding features. We can try out all feature combinations and investigate the model fit.



Note: For simplicity, we added the features in one specific (clever) order, so we cheated a bit. Also note there are  $2^{12} = 4096$  potential feature combinations.



# AVOID OVERFITTING

## Approach 4: Optimize less

Now we use polynomial regression with temperature as the only feature to predict the ozone level, i.e.,

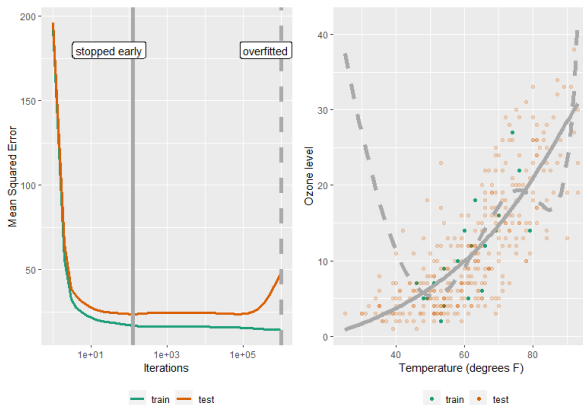
$$f(\mathbf{x} \mid \boldsymbol{\theta}) = \sum_{i=0}^d \theta_i (x_T)^i.$$

We choose  $d = 15$ , for which we get a very flexible model, which can be prone to overfitting for small data sets.

In this example, we don't solve for  $\hat{\boldsymbol{\theta}}$  directly, but instead, we use the gradient descent algorithm to find  $\hat{\boldsymbol{\theta}}$  stepwise.

# AVOID OVERFITTING

We want to stop the optimization early when the generalization error starts to degrade.



Note: For polynomial regression, gradient descent usually needs many iterations before it starts to overfit. Hence a very small training set was chosen to accelerate this effect.

# AVOID OVERFITTING

We have contradictory goals

- **maximizing the fit** (minimizing the train loss)
- **minimizing the complexity** of the model.

We need to find the “sweet spot”.



# AVOID OVERFITTING

Until now, we can either add a feature completely or not at all.

Instead of controlling the complexity in a discrete way by specifying the number of features, we might prefer to control the complexity **on a continuum** from simple to complex.



# Regularized Empirical Risk Minimization

# REGULARIZED EMPIRICAL RISK MINIMIZATION

Recall, empirical risk minimization with a complex hypothesis set tends to overfit. A major tool to handle overfitting is **regularization**.

In the broadest sense, regularization refers to any modification made to a learning algorithm that is intended to reduce its generalization error but not its training error.

Explicitly or implicitly, such modifications represent the preferences we have regarding the elements of the hypothesis set.

# REGULARIZED EMPIRICAL RISK MINIMIZATION

Commonly, regularization takes the following form:

$$\mathcal{R}_{\text{reg}}(f) = \mathcal{R}_{\text{emp}}(f) + \lambda \cdot J(f) = \sum_{i=1}^n L\left(y^{(i)}, f\left(\mathbf{x}^{(i)}\right)\right) + \lambda \cdot J(f)$$

- $J(f)$  is called **complexity penalty**, **roughness penalty** or **regularizer**.
- $\lambda > 0$  is called **complexity control** parameter.
- It measures the “complexity” of a model and penalizes it in the fit.
- As for  $\mathcal{R}_{\text{emp}}$ , often  $\mathcal{R}_{\text{reg}}$  and  $J$  are defined on  $\theta$  instead of  $f$ , so  $\mathcal{R}_{\text{reg}}(\theta) = \mathcal{R}_{\text{emp}}(\theta) + \lambda \cdot J(\theta)$ .

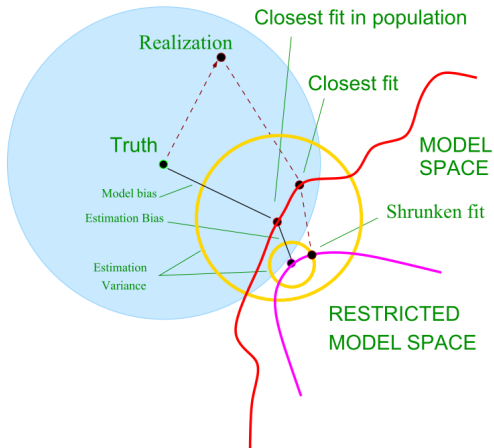
# REGULARIZED EMPIRICAL RISK MINIMIZATION

## Remarks:

- Note that we now face an optimization problem with two criteria:
  - ① models should fit well (low empirical risk),
  - ② but not be too complex (low  $J(f)$ ).
- We decide to combine the two in a weighted sum and to control the trade-off via the complexity control parameter  $\lambda$ .
- $\lambda$  is hard to set manually and is usually selected via cross-validation (see later).
- $\lambda = 0$ : The regularized risk  $\mathcal{R}_{\text{reg}}(f)$  reduces to the simple empirical  $\mathcal{R}_{\text{emp}}(f)$ .
- If  $\lambda$  goes to infinity, we stop caring about the loss/fit and models become as “simple” as possible.



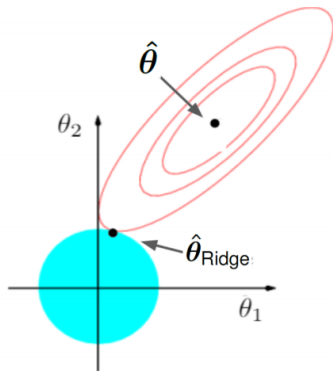
## REGULARIZED EMPIRICAL RISK MINIMIZATION



Hastie, The Elements of Statistical Learning, 2009 (p. 225)

# Einführung in das Statistische Lernen

## Lasso and Ridge Regression



### Learning goals

- Know the regularized linear model
- Know Ridge regression ( $L_2$  penalty)
- Know Lasso regression ( $L_1$  penalty)

# REGULARIZATION IN THE LINEAR MODEL

- Linear models can also overfit if we operate in a high-dimensional space with not that many observations.
- OLS usually require a full-rank design matrix.
- When features are highly correlated, the least-squares estimate becomes highly sensitive to random errors in the observed response, producing a large variance in the fit.
- We now add a complexity penalty to the loss:

$$\mathcal{R}_{\text{reg}}(\boldsymbol{\theta}) = \sum_{i=1}^n \left( y^{(i)} - \boldsymbol{\theta}^\top \mathbf{x}^{(i)} \right)^2 + \lambda \cdot J(\boldsymbol{\theta}).$$

- Intuitive to measure model complexity as deviation from the 0-origin, as the 0-model is empty and contains no effects. Models close to this either have few active features or only weak effects.
- So we measure  $J(\boldsymbol{\theta})$  through a vector norm. This shrinks coefficients closer 0, hence the term **shrinkage methods**.

# RIDGE REGRESSION

**Ridge regression** uses a simple  $L2$  penalty:

$$\begin{aligned}\hat{\boldsymbol{\theta}}_{\text{Ridge}} &= \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^n \left( y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)} \right)^2 + \lambda \|\boldsymbol{\theta}\|_2^2 \\ &= \arg \min_{\boldsymbol{\theta}} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) + \lambda \boldsymbol{\theta}^\top \boldsymbol{\theta}.\end{aligned}$$

Optimization is possible (as in the normal LM) in analytical form:

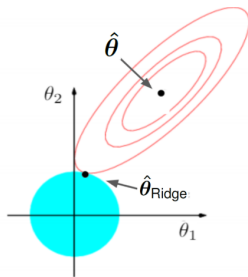
$$\hat{\boldsymbol{\theta}}_{\text{Ridge}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

Name comes from the fact that we add positive entries along the diagonal "ridge"  $\mathbf{X}^T \mathbf{X}$ .

# RIDGE REGRESSION

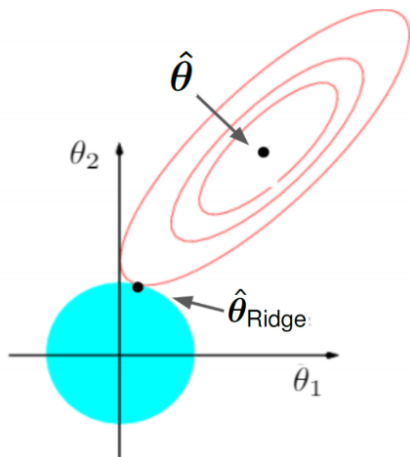
We understand the geometry of these 2 mixed components in our regularized risk objective much better, if we formulate the optimization as a constrained problem (see this a Lagrange multipliers in reverse).

$$\begin{aligned} \min_{\boldsymbol{\theta}} \quad & \sum_{i=1}^n \left( y^{(i)} - f(\mathbf{x}^{(i)} | \boldsymbol{\theta}) \right)^2 \\ \text{s.t.} \quad & \|\boldsymbol{\theta}\|_2^2 \leq t \end{aligned}$$



NB: Relationship between  $\lambda$  and  $t$  will be explained later.

# RIDGE REGRESSION



- We still optimize the  $\mathcal{R}_{emp}(\theta)$ , but cannot leave a ball around the origin.
- $\mathcal{R}_{emp}(\theta)$  grows monotonically if we move away from  $\hat{\theta}$ .
- Inside constraints perspective: From origin, jump from contour line to contour line (better) until you become infeasible, stop before.
- Outside constraints perspective: From  $\hat{\theta}$ , jump from contour line to contour line (worse) until you become feasible, stop then.
- So our new optimum will lie on the boundary of that ball.

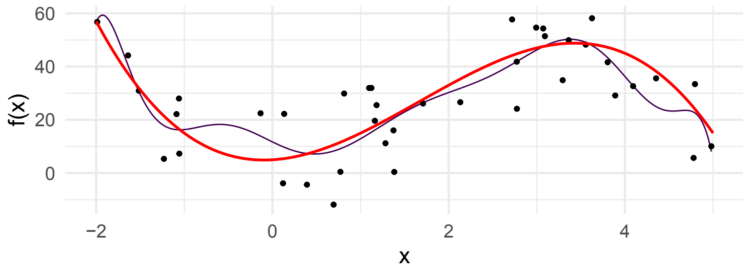
# EXAMPLE: POLYNOMIAL RIDGE REGRESSION

True (unknown) function is  $f(x) = 5 + 2x + 10x^2 - 2x^3 + \epsilon$  (in red).

Let us consider a  $d$ th-order polynomial

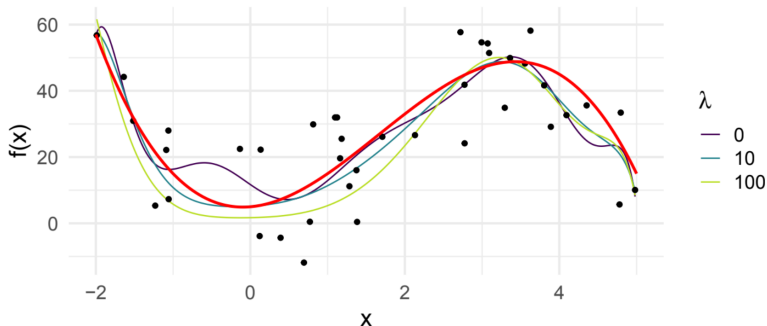
$$f(x) = \theta_0 + \theta_1 x + \dots + \theta_d x^d = \sum_{j=0}^d \theta_j x^j.$$

Using model complexity  $d = 10$  overfits:



# EXAMPLE: POLYNOMIAL RIDGE REGRESSION

With an  $L2$  penalty we can now select  $d$  "too large" but regularize our model by shrinking its coefficients. Otherwise we have to optimize over the discrete  $d$ .



$\lambda$	$\beta_0$	$\beta_1$	$\beta_2$	$\beta_3$	$\beta_4$	$\beta_5$	$\beta_6$	$\beta_7$	$\beta_8$	$\beta_9$	$\beta_{10}$
0.00	12.00	-16.00	4.80	23.00	-5.40	-9.30	4.20	0.53	-0.63	0.13	-0.01
10.00	5.20	1.30	3.70	0.69	1.90	-2.00	0.47	0.20	-0.14	0.03	-0.00
100.00	1.70	0.46	1.80	0.25	1.80	-0.94	0.34	-0.01	-0.06	0.02	-0.00



# LASSO REGRESSION

Another shrinkage method is the so-called **Lasso regression**, which uses an  $L_1$  penalty on  $\theta$ :

$$\begin{aligned}\hat{\theta}_{\text{Lasso}} &= \arg \min_{\theta} \sum_{i=1}^n \left( y^{(i)} - \theta^T \mathbf{x}^{(i)} \right)^2 + \lambda \|\theta\|_1 \\ &= \arg \min_{\theta} (\mathbf{y} - \mathbf{X}\theta)^T (\mathbf{y} - \mathbf{X}\theta) + \lambda \|\theta\|_1.\end{aligned}$$

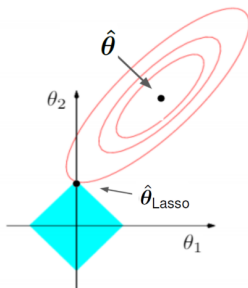
Note that optimization now becomes much harder.  $\mathcal{R}_{\text{reg}}(\theta)$  is still convex, but we have moved from an optimization problem with an analytical solution towards a non-differentiable problem.

Name: least absolute shrinkage and selection operator.

# LASSO REGRESSION

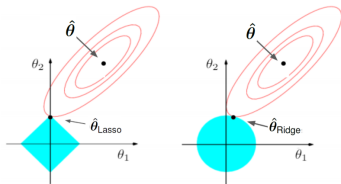
We can also rewrite this as a constrained optimization problem. The penalty results in the constrained region to look like a diamond shape.

$$\begin{aligned} \min_{\boldsymbol{\theta}} \quad & \sum_{i=1}^n \left( y^{(i)} - f(\mathbf{x}^{(i)} \mid \boldsymbol{\theta}) \right)^2 \\ \text{subject to:} \quad & \|\boldsymbol{\theta}\|_1 \leq t \end{aligned}$$



# Einführung in das Statistische Lernen

## Lasso vs. Ridge Regression

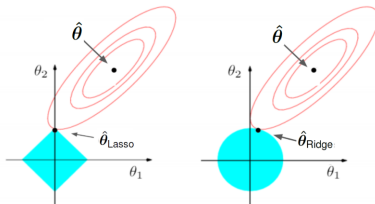


### Learning goals

- Know the geometry of Ridge vs. Lasso regularization
- Understand the effects of the methods on model coefficients
- Understand that Lasso creates sparse solutions

# LASSO VS. RIDGE GEOMETRY

$$\min_{\boldsymbol{\theta}} \sum_{i=1}^n \left( y^{(i)} - f(\mathbf{x}^{(i)} | \boldsymbol{\theta}) \right)^2 \quad \text{s.t. } \|\boldsymbol{\theta}\|_p^p \leq t$$

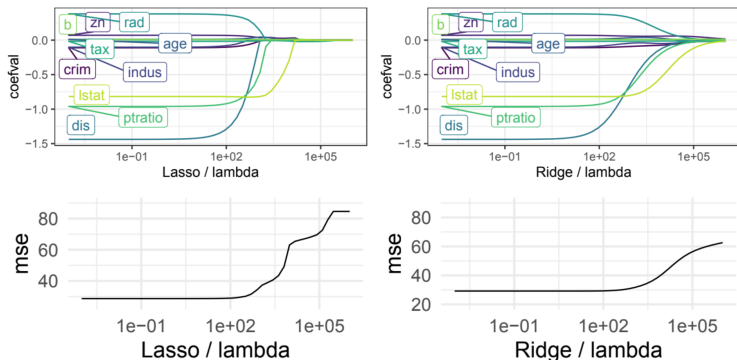


- In both cases, the solution which minimizes  $\mathcal{R}_{\text{reg}}(\boldsymbol{\theta})$  is always a point on the boundary of the feasible region (for sufficiently large  $\lambda$ ).
- As expected,  $\hat{\boldsymbol{\theta}}_{\text{Lasso}}$  and  $\hat{\boldsymbol{\theta}}_{\text{Ridge}}$  have smaller parameter norms than  $\hat{\boldsymbol{\theta}}$ .
- For Lasso, the solution likely touches vertices of the constraint region. This induces sparsity and is a form of variable selection.
- In the  $p > n$  case, the Lasso selects at most  $n$  features (due to the nature of the convex optimization problem).

# COEFFICIENT PATHS AND 0-SHRINKAGE

## Example 1: Boston Housing (few features removed for readability)

We cannot overfit here with an unregularized linear model as the task is so low-dimensional. But we see how only Lasso shrinks to sparsely 0.



Coef paths and cross-val. MSE for  $\lambda$  values for Ridge and Lasso.

# COEFFICIENT PATHS AND 0-SHRINKAGE

## Example 2: High-dimensional simulated data

We simulate a continuous, correlated dataset with 50 features, 100 observations  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(100)} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \Sigma)$  and

$$y = 10 \cdot (x_1 + x_2) + 5 \cdot (x_3 + x_4) + \sum_{j=5}^{14} x_j + \epsilon$$

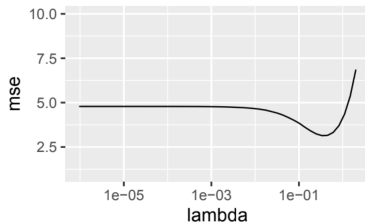
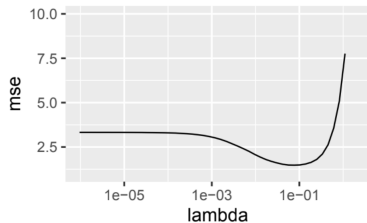
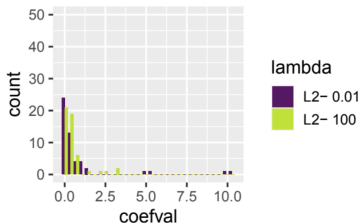
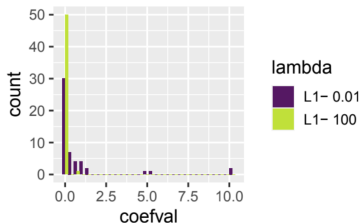
where  $\epsilon \sim \mathcal{N}(0, 1)$  and  $\forall k, l \in \{1, \dots, 50\}$ :

$$\text{Cov}(x_k, x_l) = \begin{cases} 0.7^{|k-l|} & \text{for } k \neq l \\ 1 & \text{else} \end{cases}.$$

Note that 36 of the 50 features are noise variables.

# COEFFICIENT PATHS AND 0-SHRINKAGE

Coefficient histograms for different  $\lambda$  values for Ridge and Lasso, on high-dimensional data along with the cross-validated MSE.



# REGULARIZATION AND FEATURE SCALING

- Note that very often we do not include  $\theta_0$  in the penalty term  $J(\theta)$  (but this can be implementation-dependent).
- These methods are typically not equivariant under scaling of the inputs, so one usually standardizes the features.
- Note that for a normal LM, if you scale some features, we can simply "anti-scale" the coefficients the same way. The risk does not change. For regularized models this is not so simple. If you scale features to smaller values, coefficients have to become larger to counteract. They now are penalized more heavily in  $J(\theta)$ . Such a scaling would make some features less attractive without changing anything relevant in the data.



# REGULARIZATION AND FEATURE SCALING

## Example:

- Let the true data generating process be

$$y = x_1 + \epsilon, \quad \epsilon \sim \mathcal{N}(0, 1).$$

- Let there be 5 features  $x_1, \dots, x_5 \sim \mathcal{N}(0, 1)$ .
- Using the Lasso (package `glmnet`), we get

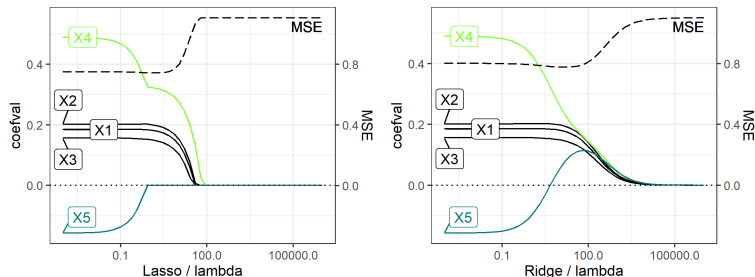
(Intercept)	x1	x2	x3	x4
-0.056	0.489	0.000	0.000	0.000

- But if we rescale any of the noise features, say  $x_2 = 10000 \cdot x_2$  and don't use standardization, we get

(Intercept)	x1	x2	x3	x4
-0.106830	0.000000	-0.000013	0.000000	0.000000

- This is due to the fact, that the coefficient of  $x_2$  will live on a very small scale as the covariate itself is large. The feature will thus get less penalized by the  $L_1$ -norm and is favored by Lasso.

# CORRELATED FEATURES



Fictional example for the model

$y = 0.2X_1 + 0.2X_2 + 0.2X_3 + 0.2X_4 + 0.2X_5 + \epsilon$  of 100 observations,  $\epsilon \sim \mathcal{N}(0, 1)$ .  $X_1$ - $X_4$  are independently drawn from different normal distributions:  $X_1, X_2, X_3, X_4 \sim \mathcal{N}(0, 2)$ . While  $X_1$ - $X_4$  have pairwise correlation coefficients of 0,  $X_4$  and  $X_5$  are nearly perfectly correlated:  $X_5 = X_4 + \delta, \delta \sim \mathcal{N}(0, 0.3), \rho(X_4, X_5) = 0.98$ .

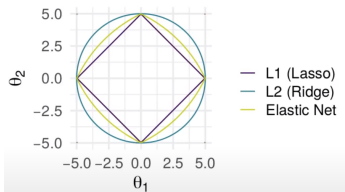
We see that Lasso shrinks the coefficient for  $X_5$  to zero early on, while Ridge assigns similar coefficients to  $X_4, X_5$  for larger  $\lambda$ .

# SUMMARIZING COMMENTS

- Neither one can be classified as overall better.
- Lasso is likely better if the true underlying structure is sparse, so if only few features influence  $y$ . Ridge works well if there are many influential features.
- Lasso can set some coefficients to zero, thus performing variable selection, while Ridge regression usually leads to smaller estimated coefficients, but still dense  $\theta$  vectors.
- Lasso has difficulties handling correlated predictors. For high correlation Ridge dominates Lasso in performance.
- For Lasso one of the correlated predictors will have a larger coefficient, while the rest are (nearly) zeroed. The respective feature is, however, selected randomly.
- For Ridge the coefficients of correlated features are similar.

# Einführung in das Statistische Lernen

## Elastic Net and Regularization for GLMs



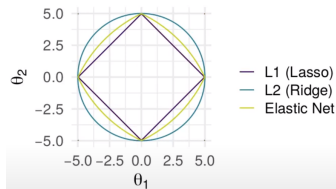
### Learning goals

- Know the elastic net as compromise between Ridge and Lasso regression
- Know regularized logistic regression

# ELASTIC NET

Elastic Net combines the  $L_1$  and  $L_2$  penalties:

$$\mathcal{R}_{\text{elnet}}(\boldsymbol{\theta}) = \sum_{i=1}^n (y^{(i)} - \boldsymbol{\theta}^\top \mathbf{x}^{(i)})^2 + \lambda_1 \|\boldsymbol{\theta}\|_1 + \lambda_2 \|\boldsymbol{\theta}\|_2^2.$$



- Correlated predictors tend to be either selected or zeroed out together.
- Selection of more than  $n$  features possible for  $p > n$ .

# ELASTIC NET

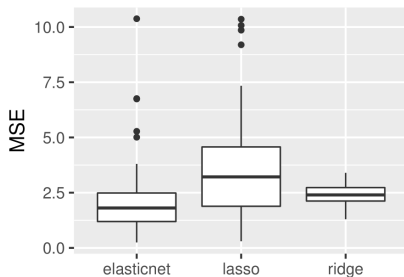
Simulating two examples with each 50 data sets and 100 observations each:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \sigma\epsilon, \quad \epsilon \sim N(0, 1), \quad \sigma = 1$$

**Ridge** performs better for:

$$\boldsymbol{\beta} = (\underbrace{2, \dots, 2}_5, \underbrace{0, \dots, 0}_5)$$

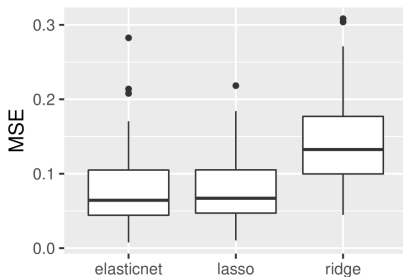
$$\text{corr}(\mathbf{X}_i, \mathbf{X}_j) = 0.8^{|i-j|} \text{ for all } i \text{ and } j$$



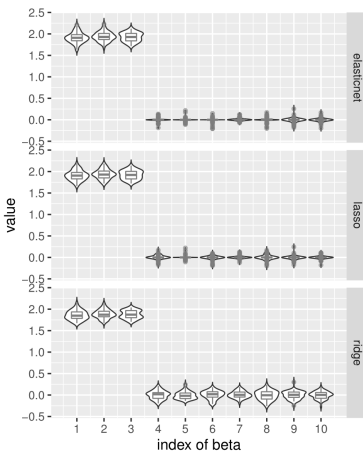
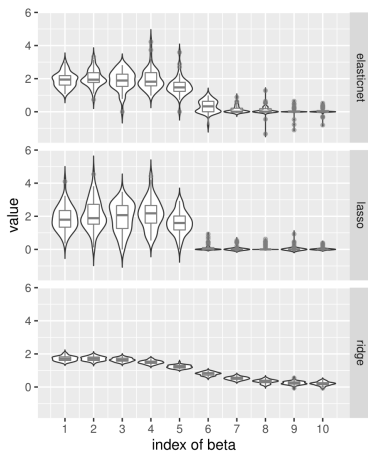
**Lasso** performs better for:

$$\boldsymbol{\beta} = (2, 2, 2, \underbrace{0, \dots, 0}_7)$$

$$\text{corr}(\mathbf{X}_i, \mathbf{X}_j) = 0 \text{ for all } i \neq j, \text{ otherwise } 1$$



# ELASTIC NET



Since Elastic Net offers a compromise between Ridge and Lasso, it is suitable for both data situations.

# REGULARIZED LOGISTIC REGRESSION

Regularizers can be added very flexibly to basically any model which is based on ERM.

Hence, we can, e.g., construct  $L_1$ - or  $L_2$ -penalized logistic regression to enable coefficient shrinkage and variable selection in this model.

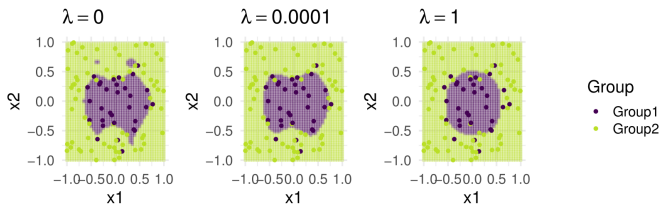
$$\begin{aligned}\mathcal{R}_{\text{reg}}(\boldsymbol{\theta}) &= \mathcal{R}_{\text{emp}}(\boldsymbol{\theta}) + \lambda \cdot J(\boldsymbol{\theta}) \\ &= \sum_{i=1}^n \log \left[ 1 + \exp \left( -2y^{(i)} f \left( \mathbf{x}^{(i)} \mid \boldsymbol{\theta} \right) \right) \right] + \lambda \cdot J(\boldsymbol{\theta})\end{aligned}$$



# REGULARIZED LOGISTIC REGRESSION

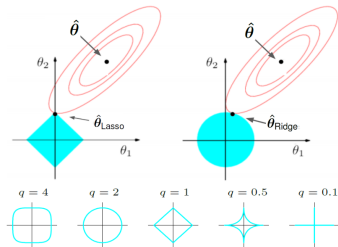
We fit a logistic regression model using polynomial features for  $x_1$  and  $x_2$  with maximum degree of 7. We add an  $L_2$  penalty. We see for

- $\lambda = 0$ : The unregularized model seems to overfit.
- $\lambda = 0.0001$ : Regularization helps to learn the underlying mechanism.
- $\lambda = 1$ : The real data-generating process is captured very well.



# Einführung in das Statistische Lernen

## L0 Regularization

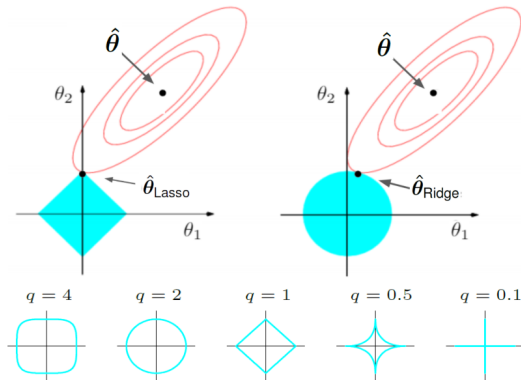


### Learning goals

- Know LQ norm regularization
- Understand that L0 norm realization simply counts the number of non-zero parameters

# LQ NORM REGULARIZATION

Besides  $L_1$  and  $L_2$  norm we could use any  $L_q$  norm for regularization.



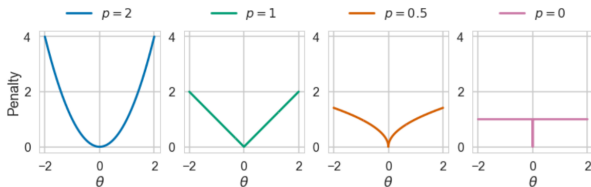
**Figure:** *Top:* Ridge and Lasso loss contours and feasible regions. *Bottom:* Different feasible region shapes for  $L_q$  norms  $\sum_j |\theta_j|^q$ .

# L0 REGULARIZATION

- Consider the  $L_0$ -regularized risk of a model  $f(\mathbf{x} \mid \boldsymbol{\theta})$

$$\mathcal{R}_{\text{reg}}(\boldsymbol{\theta}) = \mathcal{R}_{\text{emp}}(\boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}\|_0 := \mathcal{R}_{\text{emp}}(\boldsymbol{\theta}) + \lambda \sum_j |\theta_j|^0.$$

- Unlike the  $L_1$  and  $L_2$  norms, the  $L_0$  "norm" simply counts the number of non-zero parameters in the model.



Credit: Christos Louizos

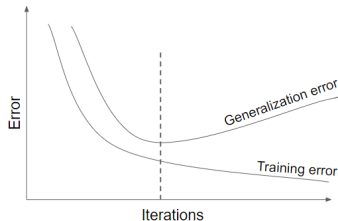
**Figure:**  $L_p$  norm penalties for a parameter  $\boldsymbol{\theta}$  according to different values of  $p$ .

# L0 REGULARIZATION

- For any parameter  $\theta$ , the  $L_0$  penalty is zero for  $\theta = 0$  (defining  $0^0 := 0$ ) and is constant for any  $\theta \neq 0$ , no matter how large or small it is.
- $L_0$  regularization induces sparsity in the parameter vector more aggressively than  $L_1$  regularization, but does not shrink concrete parameter values as  $L_1$  and  $L_2$  does.
- Model selection criteria such as Akaike Information Criterion (AIC) and Bayesian Information Criterion (BIC) are special cases of  $L_0$  regularization (corresponding to specific values of  $\lambda$ ).
- The  $L_0$ -regularized risk is neither continuous, differentiable or convex.
- It is computationally hard to optimize (NP-hard) and likely intractable. For smaller  $n$  and  $p$  we might be able to solve this nowadays directly, for larger scenarios efficient approximations of the  $L_0$  are still topic of current research.

# Einführung in das Statistische Lernen

## Early Stopping

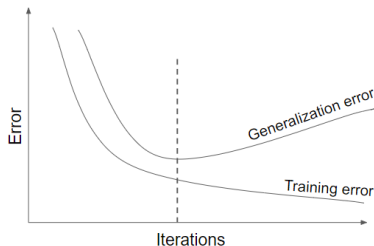


### Learning goals

- Know how early stopping works
- Understand how early stopping acts as a regularizer

# EARLY STOPPING

- When training with an iterative optimizer such as SGD, it is commonly the case that, after a certain number of iterations, generalization error begins to increase even though training error continues to decrease.
- **Early stopping** refers to stopping the algorithm early before the generalization error increases.



**Figure:** After a certain number of iterations, the algorithm begins to overfit.

# EARLY STOPPING

How early stopping works:

- ➊ Split training data  $\mathcal{D}_{\text{train}}$  into  $\mathcal{D}_{\text{subtrain}}$  and  $\mathcal{D}_{\text{val}}$  (e.g. with a ratio of 2:1).
- ➋ Train on  $\mathcal{D}_{\text{subtrain}}$  and evaluate model using the validation set  $\mathcal{D}_{\text{val}}$ .
- ➌ Stop training when validation error stops decreasing (after a range of “patience” steps).
- ➍ Use parameters of the previous step for the actual model.

More sophisticated forms also apply cross-validation.



# EARLY STOPPING

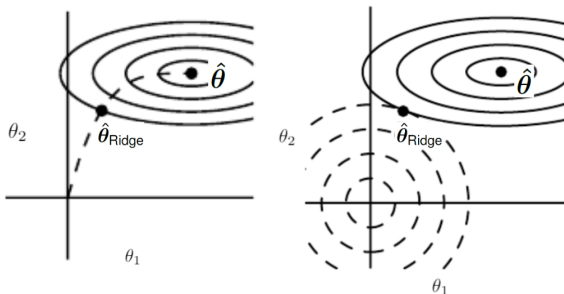
Strengths	Weaknesses
Effective and simple	Periodical evaluation of validation error
Applicable to almost any model without adjustment	Temporary copy of $\theta$ (we have to save the whole model each time validation error improves)
Combinable with other regularization methods	Less data for training $\rightarrow$ include $\mathcal{D}_{\text{val}}$ afterwards

- Relation between optimal early-stopping iteration  $T_{\text{stop}}$  and weight-decay penalization parameter  $\lambda$  for step-size  $\alpha$  (see Goodfellow et al. (2016) page 251-252 for proof):

$$T_{\text{stop}} \approx \frac{1}{\alpha\lambda} \Leftrightarrow \lambda \approx \frac{1}{T_{\text{stop}}\alpha}$$

- Small  $\lambda$  (low penalization)  $\Rightarrow$  high  $T_{\text{stop}}$  (complex model / lots of updates).

# EARLY STOPPING



Credit: Goodfellow et al. (2016)

**Figure:** An illustration of the effect of early stopping. *Left:* The solid contour lines indicate the contours of the negative log-likelihood. The dashed line indicates the trajectory taken by SGD beginning from the origin. Rather than stopping at the point  $\hat{\theta}$  that minimizes the risk, early stopping results in the trajectory stopping at an earlier point  $\hat{\theta}_{\text{Ridge}}$ . *Right:* An illustration of the effect of  $L_2$  regularization for comparison. The dashed circles indicate the contours of the  $L_2$  penalty which causes the minimum of the total cost to lie closer to the origin than the minimum of the unregularized cost.