

Machine Learning

A. Supervised Learning: Linear Models & Fundamentals A.3. Regularization

Lars Schmidt-Thieme

Information Systems and Machine Learning Lab (ISMLL) Institute for Computer Science University of Hildesheim, Germany

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning

Syllabus

Fri. 2	27.10.	(1)	0. Introduction	
			A. Supervised Learning: Linear Models & Fundamentals	
Fri.	3.11.	(2)	A.1 Linear Regression	
Fri. 1	.0.11.	(3)	A.2 Linear Classification	
Fri. 1	7.11.	(4)	A.3 Regularization	
Fri. 2	24.11.	(5)	A.4 High-dimensional Data	
B. Supervised Learning: Nonlinear Models				
Fri.	1.12.	(6)	B.1 Nearest-Neighbor Models	
Fri.	8.12.	(7)	B.2 Neural Networks	
Fri. 1	.5.12.	(8)	B.3 Decision Trees	
Fri.	12.1.	(9)	B.4 Support Vector Machines	
Fri.	19.1.	(10)	B.5 A First Look at Bayesian and Markov Networks	
			C. Unsupervised Learning	
Fri.	26.1.	(11)	C.1 Clustering	
Fri.	2.2.	(12)	C.2 Dimensionality Reduction	
Fri.	9.2.	(13)	C.3 Frequent Pattern Mining	



Outline



- 1. The Problem of Overfitting
- 2. Model Selection
- 3. Regularization
- 4. Hyperparameter Optimization

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning

Outline

1. The Problem of Overfitting

- 2. Model Selection
- 3. Regularization
- 4. Hyperparameter Optimization









Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning

${\sf Underfitting}/{\sf Overfitting}$

Underfitting:

- the model is not complex enough to explain the data well.
- ► results in poor predictive performance.

Overfitting:

- the model is too complex, it describes the
 - noise instead of the
 - underlying relationship between target and predictors.
- results in poor predictive performance as well.

Remark: Given N points (x_n, y_n) without repeated measurements (i.e. $x_n \neq x_m, n \neq m$), there exists a polynomial of degree N - 1 with RSS equal to 0.



Outline



1. The Problem of Overfitting

2. Model Selection

- 3. Regularization
- 4. Hyperparameter Optimization

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany



Machine Learning

Losses and Fit Measures

	loss	fit/quality measure
semantics	the smaller, the better	the larger, the better
goal	minimize	maximize
	$RSS(y, \hat{y})$	$\log L_{\mathcal{N}}(y, \hat{y})$
regression	$= \sum_{n=1}^{N} (y_n - \hat{y}_n)^2$	$= \sum_{n=1}^{N} - \frac{1}{2\sigma_{y}^{2}} (y_{n} - \hat{y}_{n})^{2}$
	$RMSE(y, \hat{y})$	
	$:= (\frac{1}{N} \sum_{n=1}^{N} (y_n - \hat{y}_n)^2)^{\frac{1}{2}}$	
	$MAE(y, \hat{y})$	
	$:= \frac{1}{N} \sum_{n=1}^{N} y_n - \hat{y}_n $	
classification	$MR(y, \hat{y})$	$ACC(y, \hat{y})$
classification	$:=\sum_{n=1}^{N}\mathbb{I}(y_n\neq \hat{y}_n)$	$:=\sum_{n=1}^{N}\mathbb{I}(y_n=\hat{y}_n)$
		$\log L_{binomial}(y, \hat{y})$
		$:=\sum_{n=1}^{N}\pi\mathbb{I}(y_n=\hat{y}_n)$
		$+(1-\pi)\mathbb{I}(y_n eq \hat{y}_n)$

Model Selection Measures

Model selection: given a set of models, e.g.,

$$Y = \sum_{m=0}^{p-1} \beta_m X_m$$

indexed by p (i.e., one model for each value of p), make a choice which model **describes** the data best.

If we just look at losses / fit measures such as RSS, then

the larger p, the better the fit

or equivalently

the larger p, the lower the loss

as the model with p parameters can be **reparametrized** in a model with p' > p parameters by setting

$$eta'_m = \left\{ egin{array}{cc} eta_m, & ext{for } m \leq p \\ 0, & ext{for } m > p \end{array}
ight.$$

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning

Model Selection Measures

```
One uses model selection measures of type
```

```
model selection measure = fit – complexity (max!)
```

or equivalently

model selection measure = loss + complexity (min!)

The smaller the loss (= lack of fit), the better the model.

The smaller the complexity, the simpler and thus better the model.

The model selection measure tries to find a trade-off between fit/loss and complexity.





Model Selection Measures



Akaike Information Criterion (AIC):

(maximize)

AIC := $\log L - p$

or (minimize)

$$AIC := -2\log L + 2p$$

Bayes Information Criterion (BIC) / Bayes-Schwarz Information Criterion: (maximize)

$$\mathsf{BIC} := \log L - \frac{p}{2} \log N$$

where L denotes the likelihood, N the number of samples.

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning

Variable Backward Selection

{ A, F, H, I, J, L, P } AIC = 63.01



Variable Backward Selection





Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning



Variable Backward Selection



Variable Backward Selection





X removed variable

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning

Outline

- 1. The Problem of Overfitting
- 2. Model Selection
- 3. Regularization
- 4. Hyperparameter Optimization

7 / 25

Shrinkage



Model selection operates by

- 1. fitting model instances for a set of models with varying complexity
- 2. picking the "best one" ex post,

Variable Selection

- ▶ = model selection applied to models with different predictor subsets
- for models ŷ that factor through a linear combination of the predictors,

$$\hat{y}(x;\hat{eta})=f(\sum_{m=1}^{M}\hat{eta}_mx_m)$$
 for a suitable f

- dropping a variable x_m from the model is equivalent to
- forcing its model parameter $\hat{\beta}_m$ to be 0.

Note: "Fitting a model instance" = "Learning model parameters", for models having parameters such as linear regression, logistic regression etc. Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning

Shrinkage

Variable Selection

- ▶ ...
 - forcing its model parameter $\hat{\beta}_m$ to be 0.

Shrinkage follows a similar idea:

- ► smaller parameters mean a simpler hypothesis/less complex model.
- hence, small parameters should be prefered in general.
- a term is added to the objective function to
 - favor small parameters or equivalently
 - penalize large parameters or
 - shrink them towards 0

instead of forcing them to be 0.



Shrinkage



There are various types of shrinkage techniques for different application domains.

L1/Lasso Regularization: $\lambda \sum_{m=1}^{M} \left| \hat{\beta}_{m} \right| = \lambda \left\| \hat{\beta} \right\|_{1}$

L2/Tikhonov Regularization: $\lambda \sum_{m=1}^{M} \hat{\beta}_{m}^{2} = \lambda \left\| \hat{\beta} \right\|_{2}^{2}$

Elastic Net: $\lambda_1 \left\| \hat{\beta} \right\|_1 + \lambda_2 \left\| \hat{\beta} \right\|_2^2$

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning

Ridge Regression



Ridge regression is a combination of



= L2 loss $+\lambda$ L2 regularization

Ridge Regression (Closed Form) Ridge regression: minimize

$$\mathsf{RSS}_{\lambda}(\hat{\beta}) = \mathsf{RSS}(\hat{\beta}) + \lambda \sum_{j=1}^{p} \hat{\beta}_{j}^{2} = \langle \mathbf{y} - \mathbf{X}\hat{\beta}, \mathbf{y} - \mathbf{X}\hat{\beta} \rangle + \lambda \sum_{j=1}^{p} \hat{\beta}_{j}^{2}$$
$$\Rightarrow \hat{\beta} = \left(\mathbf{X}^{T}\mathbf{X} + 2\lambda I\right)^{-1} \mathbf{X}^{T}\mathbf{y}, \quad I := \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & 1 \end{pmatrix}$$

with $\lambda \geq 0$ a complexity parameter / regularization parameter.

Beware: ridge regression parameter estimates are not equivariant under scaling of the predictors

→ data should be normalized before parameter estimation:

$$x'_{n,m} := \frac{x_{n,m} - \bar{x}_{.,m}}{\hat{\sigma}(x_{.,m})}$$

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning

Ridge Regression (Gradient Descent)

1 **learn-ridgereg-GD**($\mathcal{D}^{train} := \{(x_1, y_1), \dots, (x_N, y_N)\}, \alpha, t_{max} \in \mathbb{N}, \epsilon \in \mathbb{R}^+\}$:

2
$$X := (x_1, x_2, ..., x_N)^T$$

3 $y := (y_1, y_2, ..., y_N)^T$
4 $\hat{\beta} := 0_M$
5 $\ell := ||y - X\hat{\beta}||^2$

6 for
$$t = 1, ..., t_{max}$$
:

7
$$\hat{\beta} := \hat{\beta} - \alpha (-2 \cdot X^T (y - X\hat{\beta}) + 2\lambda\hat{\beta})$$

$$_{8}$$
 $\ell^{\mathsf{old}} := \ell$

9
$$\ell \coloneqq ||y - X\hat{eta}||^2$$

10 if
$$\ell - \ell^{\mathsf{old}} < \epsilon$$
:

11 return $\hat{\beta}$

raise exception "not converged in t_{max} iterations"

L2-Regularized Update Rule

$$\hat{\beta}^{(t)} := \underbrace{(1 - 2\alpha\lambda)}_{\text{shrinkage}} \hat{\beta}^{(t-1)} - \alpha \left(-2X^{T} (y - X\hat{\beta}^{(t-1)}) \right)$$







Tikhonov Regularization Derivation (1/2)

Treat the true parameters θ_j as random variables Θ_j with the following distribution (prior):

$$\Theta_j \sim \mathcal{N}(0, \sigma_{\Theta}), \quad j = 1, \dots, p$$

Then the joint likelihood of the data and the parameters is

$$L_{\mathcal{D},\Theta}(\theta) := \left(\prod_{n=1}^{N} p(x_n, y_n \mid \theta)\right) \prod_{j=1}^{p} p(\Theta_j = \theta_j)$$

and the conditional joint log likelihood of the data and the parameters

$$\log L_{\mathcal{D},\Theta}^{\text{cond}}(\theta) := \left(\sum_{n=1}^{N} \log p(y_n \,|\, x_n, \theta)\right) + \sum_{j=1}^{p} \log p(\Theta_j = \theta_j)$$

and

$$\log p(\Theta_j = \theta_j) = \log \frac{1}{\sqrt{2\pi}\sigma_{\Theta}} e^{-\frac{\theta_j^2}{2\sigma_{\Theta}^2}} = -\log(\sqrt{2\pi}\sigma_{\Theta}) - \frac{\theta_j^2}{2\sigma_{\Theta}^2}$$

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning

Tikhonov Regularization Derivation (2/2)Dropping the terms that do not depend on θ_j yields:

$$\log L_{\mathcal{D},\Theta}^{\text{cond}}(\theta) := \left(\sum_{n=1}^{N} \log p(y_n | x_n, \theta)\right) + \sum_{j=1}^{p} \log p(\Theta_j = \theta_j)$$
$$\propto \left(\sum_{n=1}^{N} \log p(y_n | x_n, \theta)\right) - \frac{1}{2\sigma_{\Theta}^2} \sum_{j=1}^{p} \theta_j^2$$

This also gives a semantics to the complexity / regularization parameter λ :

$$\lambda = \frac{1}{2\sigma_{\Theta}^2}$$

but σ_{Θ}^2 is unknown. (We will see methods to estimate λ soon.)

The parameters θ that maximize the joint likelihood of the data and the parameters are called Maximum Aposteriori Estimators (MAP estimators).





L2-Regularized Logistic Regression (Gradient Descent)



$$\log L_{\mathcal{D}}^{\text{cond}}(\hat{\beta}) = \sum_{n=1}^{N} y_n \langle x_n, \hat{\beta} \rangle - \log(1 + e^{\langle x_n, \hat{\beta} \rangle}) - 2\lambda \sum_{j=1}^{P} \hat{\beta}_j^2$$

1: procedure LOG-REGR-
GA(
$$\mathcal{L}_{\mathcal{D}}^{\text{cond}}$$
 : $\mathbb{R}^{P+1} \to \mathbb{R}, \hat{\beta}^{(0)} \in \mathbb{R}^{P+1}, \alpha, t_{\max} \in \mathbb{N}, \epsilon \in \mathbb{R}^{+}$)
2: for $t = 1, \ldots, t_{\max}$ do
3: $\hat{\beta}_{0}^{(t)} := \hat{\beta}_{0}^{(t-1)} + \alpha \sum_{n=1}^{N} \left(y_{n} - p\left(Y = 1 | X = x_{n}; \hat{\beta}^{(t-1)}\right) \right)$
4: for $j = 1, \ldots, P$ do
5: $\hat{\beta}_{j}^{(t)} :=$
 $\hat{\beta}_{j}^{(t-1)} + \alpha \left(\sum_{n=1}^{N} x_{n,j} \left(y_{n} - p\left(Y = 1 | X = x_{i}; \hat{\beta}^{(t-1)}\right) \right) - 2\lambda \hat{\beta}_{j}^{(t-1)} \right)$

6: if
$$L_{\mathcal{D}}^{\text{cond}}(\hat{\beta}^{(t-1)}) - L_{\mathcal{D}}^{\text{cond}}(\hat{\beta}^{(t)})) < \epsilon$$
 then
7: return $\hat{\beta}^{(t)}$

8: **error** "not converged in
$$t_{max}$$
 iterations"

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning

L2-Regularized Logistic Regression (Newton)

Newton update rule:

$$\hat{\beta}^{(t)} := \hat{\beta}^{(t-1)} + \alpha H^{-1} \nabla_{\hat{\beta}} p\left(Y = 1 | X = x_i; \hat{\beta}^{(t-1)}\right)$$

$$p_i = p\left(Y = 1 | X = x_i; \hat{\beta}^{(t-1)}\right)$$

$$\nabla_{\hat{\beta}} L_{\mathcal{D}}^{cond} = \begin{pmatrix} \sum_{n=1}^{N} -(y_n - p_n) \\ \sum_{n=1}^{N} -x_{n,1} (y_n - p_n) - 2\lambda \hat{\beta}_1 \\ \vdots \\ \sum_{n=1}^{N} -x_{n,P} (y_n - p_n) - 2\lambda \hat{\beta}_P \end{pmatrix}$$

$$H = \sum_{n=1}^{N} -p_n (1 - p_n) x_n x_n^T - 2\lambda I$$



Outline



- 1. The Problem of Overfitting
- 2. Model Selection
- 3. Regularization

4. Hyperparameter Optimization

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning



- Most models and learning algorithms have parameters that cannot be learned by minimizing the objective function, because either
 - the objective function would be minimized for a trivial value, e.g., $\lambda=$ 0, or
 - ► the parameters affect the learning algorithm, e.g., learning rate.
- These parameters are called hyperparameters λ and they parametrize a learning algorithm A_{λ} .
 - choose suitable hyperparameters λ
 - ► use A_λ to map the training data D_{train} to a prediction function ŷ by minimizing some loss L(D, ŷ) over the training data.

What is Hyperparameter Optimization?

- Identifying good values for the hyperparameters λ is called hyperparameter optimization.
 - hyperparameter optimization is a second level optimization

$$\argmin_{\lambda \in \Lambda} \mathcal{L}(\mathcal{D}_{\mathsf{valid}}, \mathcal{A}_{\lambda}(\mathcal{D}_{\mathsf{train}})) = \argmin_{\lambda \in \Lambda} \Psi(\lambda)$$

where

- Ψ is the hyperparameter response function and
- D_{valid} a validation data
 (aka calibration data and holdout data).

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning

Why Hyperparameter Optimization

- ► So far only model parameters were optimized.
- ► Values for hyperparameters (such as regularization λ and learning rate α) came "out of the blue".
- Hyperparameters can have a big impact on the prediction quality.









Jniversizar

Grid Search

- Assume we have Q hyperparameters $\lambda_1, \ldots, \lambda_Q$
- Choose for each hyperparameter λ_q a set of values Λ_q .
- $\Lambda := \prod_{q=1}^{Q} \Lambda_q$ is then a grid of hyperparameters.
- Choose the hyperparameter combination $\lambda \in \Lambda$ with best performance on \mathcal{D}_{valid} .



Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning

Random Search

- Instead of trying hyperparameter combinations on a grid, try random hyperparameter combinations λ for Λ (within a reasonable range).
- Usually slightly better results than grid search.









- - Whenever a learning process depends on a hyperparameter, the hyperparameter can be estimated by picking the value with the lowest error.
 - If this is done on test data, one actually uses test data in the training process ("train on test"), thereby lessen its usefulness for estimating the test error.
 - Therefore, one splits the training data again in
 - (proper) training data and
 - validation data.
 - The validation data figures as test data during the training process.

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning

Cross Validation



Instead of a single split into

training data, (validation data,) and test data

K-fold cross validation splits the data in K parts (of roughly equal size)

$$\mathcal{D} = \mathcal{D}_1 \cup \mathcal{D}_2 \cup \dots \cup \mathcal{D}_{\mathcal{K}}, \quad \mathcal{D}_k$$
 pairwise disjoint

and averages performance over K learning problems

 $\mathcal{D}_{\mathsf{train}}^{(k)} \coloneqq \mathcal{D} \setminus \mathcal{D}_k, \quad \mathcal{D}_{\mathsf{test}}^{(k)} \coloneqq \mathcal{D}_k \quad k = 1, \dots, K$

Common is 5- and 10-fold cross validation.

N-fold cross validation is also known as leave one out.



Cross Validation

How many folds to use in K-fold cross validation?

K = N / leave one out:

- approximately unbiased for the true prediction error.
- high variance as the N training sets are very similar.
- in general computationally costly as N different models have to be learnt.

K = 5:

- ► lower variance.
- bias could be a problem, due to smaller training set size the prediction error could be overestimated.

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning

Summary

- The problem of underfitting can be overcome by using more complex models, e.g., having
 - variable interactions as in polynomial models.
- The problem of overfitting can be overcome by
 - model selection / variable selection as well as by
 - (parameter) shrinkage.
- Applying L2-regularization to Linear and Logistic Regression requires only few changes in the learning algorithm
- Shrinkage introduces a hyperparameter λ that cannot be learned by direct loss minimization.
- Estimating the best hyperparameters can be considered as a meta-learning problem. They can be estimated e.g. by
 - Grid Search or
 - Random Search.

using validation data.



Further Readings



[James et al., 2013, chapter 3], [Murphy, 2012, chapter 7], [Hastie et al., 2005, chapter 3].

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning

References



Trevor Hastie, Robert Tibshirani, Jerome Friedman, and James Franklin. The Elements of Statistical Learning: Data Mining, Inference and Prediction, volume 27. Springer, 2005.

Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *An Introduction to Statistical Learning*. Springer, 2013. Kevin P. Murphy. *Machine Learning: A Probabilistic Perspective*. The MIT Press, 2012.