

Regularization and model selection

Chapter 4

1

Reminder: bias vs variance, overfitting

\mathcal{H}_0 \mathcal{H}_1

bias = 0.50 var = 0.25 bias = 0.21 var = 1.69

2

Regularization

The cure for overfitting - regularization

Without regularization With regularization

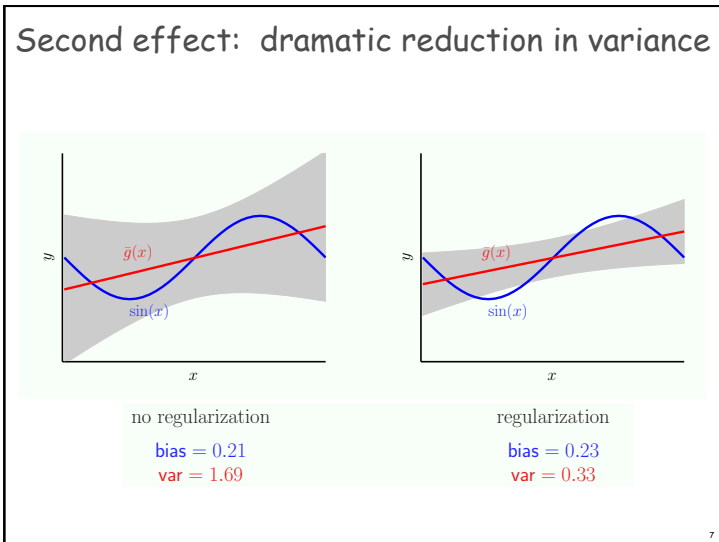
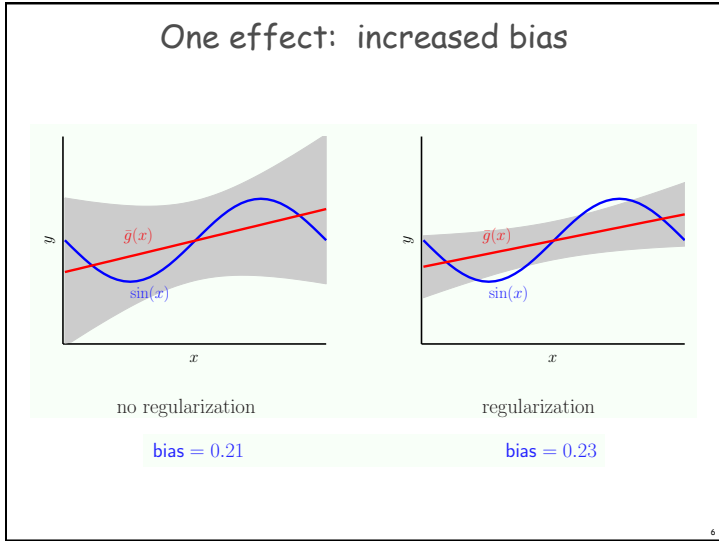
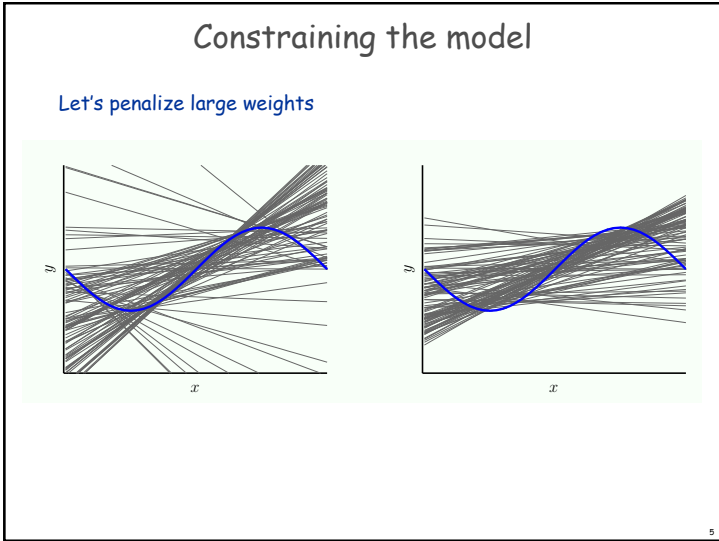
3

Regularization

How does it work?

- ❖ **Constrains** the model so it cannot fit the noise
- ❖ Potential side effect: if it cannot fit the noise, can it fit the target function?
- ❖ **Introduces bias and reduces variance**, so that (hopefully) out-of-sample error is lower

4



Constraining the complexity of the model

Replace E_{in} with:

$$E_{aug}(h) = E_{in}(h) + \frac{\lambda}{N} \Omega(h)$$

λ regularization constant Regularization term

E_{aug} is a better proxy for E_{out} than E_{in}

8

Choosing a regularizer

We want to constrain the learned function in the direction of the target function.

Intuition: noise is non-smooth

Common choice for the augmented in-sample-error:

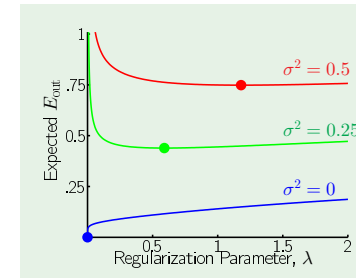
$$E_{aug}(\mathbf{w}) = E_{in}(\mathbf{w}) + \lambda \mathbf{w}^T \mathbf{w}$$

weight decay regularizer

9

Is there an optimal value for λ ?

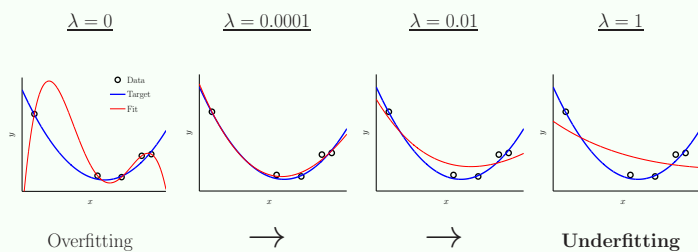
The behavior of E_{out} as a function of the regularization parameter for varying levels of noise:



10

Is there an optimal value for λ ?

Minimizing $E_{aug}(\mathbf{w}) = E_{in}(\mathbf{w}) + \lambda \mathbf{w}^T \mathbf{w}$



11

Regularized least-squares

Ridge regression:

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) + \lambda \|\mathbf{w}\|^2$$

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

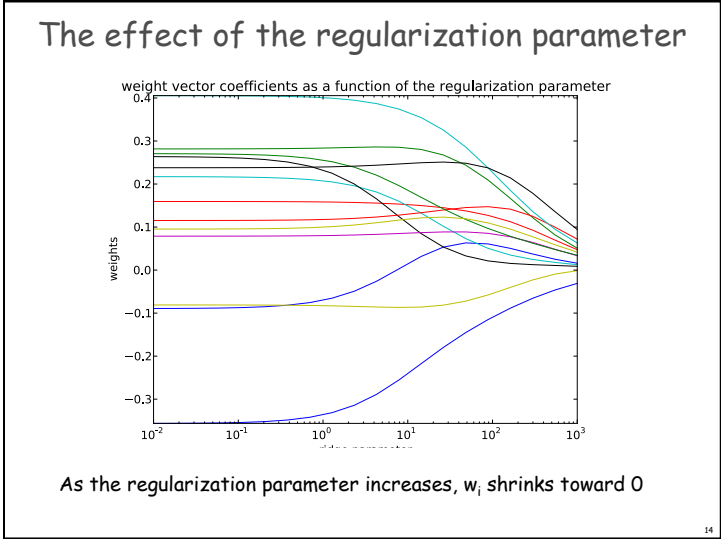
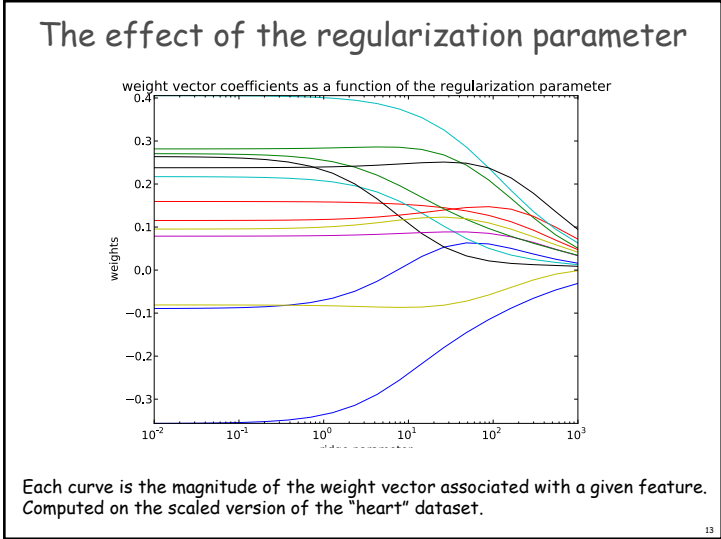
The regularization term controls the size of the components of the weight vector.

There is a tradeoff between fitting (the error term) and regularization. The regularization terms can therefore **prevent overfitting**. The parameter λ controls this tradeoff.

Many ML methods can be expressed as solution to a criterion of the form:

error term + regularization term

12



Assignment 3

Explore the effect of regularization with least-squares regression.

The validation set

How to choose the value of the regularization parameter?
 Take a sneak peak at E_{out} using a validation set!

On a validation set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_K, y_K)$, the error is $E_{\text{val}}(h) = \frac{1}{K} \sum_{k=1}^K e(h(\mathbf{x}_k), y_k)$

$$\mathbb{E}[E_{\text{val}}(h)] = \frac{1}{K} \sum_{k=1}^K \mathbb{E}[e(h(\mathbf{x}_k), y_k)] = E_{\text{out}}(h)$$

$$\text{var}[E_{\text{val}}(h)] = \frac{1}{K^2} \sum_{k=1}^K \text{var}[e(h(\mathbf{x}_k), y_k)] = \frac{\sigma^2}{K}$$

$$E_{\text{val}}(h) = E_{\text{out}}(h) \pm O\left(\frac{1}{\sqrt{K}}\right)$$

Section 4.3.1

Choosing the size of the validation set

Given the data set $\mathcal{D} = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$

$\underbrace{K \text{ points}}_{\mathcal{D}_{\text{val}}} \rightarrow \text{validation}$ $\underbrace{N - K \text{ points}}_{\mathcal{D}_{\text{train}}} \rightarrow \text{training}$

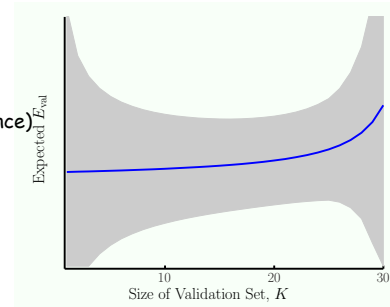
$O\left(\frac{1}{\sqrt{K}}\right)$: Small $K \implies$ bad estimate
 Large $K \implies$?

Rule of thumb: use 20% of the data for validation

17

Choosing the size of the validation set

Shaded region:
the uncertainty (variance)
of the estimate



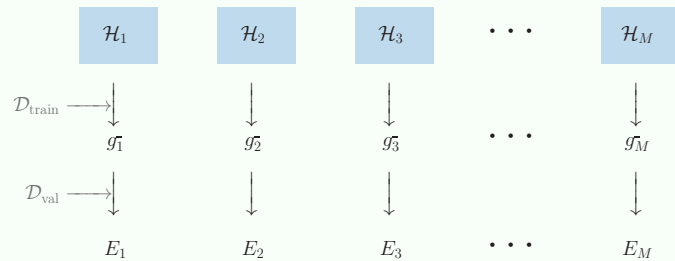
Observations:

- As we increase the size of the validation set, the estimate goes up because of a small training set
- The uncertainty in E_{val} decreases as we increase K , up to a point, where a small training set size generates uncertainty in the estimate

18

Using the validation set

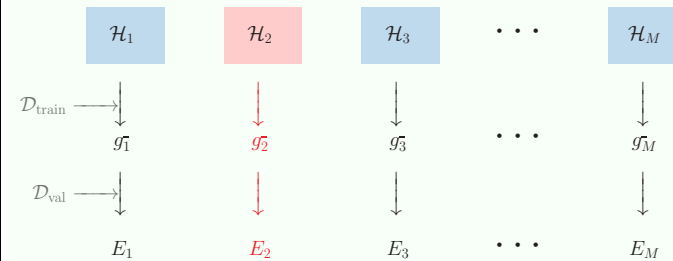
The validation set is used to get estimates that allow us to choose a value for the regularization parameter.



19

Using the validation set

The validation set is used to get estimates that allow us to choose a value for the regularization parameter.



20

Using the validation set

The validation set is used to get estimates that allow us to choose a value for the regularization parameter.

(\mathcal{H}, λ_1)

↓

g_1

(\mathcal{H}, λ_2)

↓

g_2

(\mathcal{H}, λ_3)

↓

g_3

...

(\mathcal{H}, λ_M)

↓

g_M

21

Using the validation set

M models $\mathcal{H}_1, \dots, \mathcal{H}_M$

Use $\mathcal{D}_{\text{train}}$ to learn g_m^- for each model

Evaluate g_m^- using \mathcal{D}_{val} :

$$E_m = E_{\text{val}}(g_m^-); \quad m = 1, \dots, M$$

Pick model $m = m^*$ with smallest E_m

At the end: train a model on all the data using the parameters of \mathcal{H}_{m^*} .

22

We have a dilemma...

We would like to have the following:

$$E_{\text{out}}(g) \approx E_{\text{out}}(g^-) \approx E_{\text{val}}(g^-)$$

(small K) (large K)

g : the model as a result of training on all the data
 g^- : the model trained on $\mathcal{D}_{\text{train}}$

Can we have K both large and small?

23

Leave-one-out errors

Extreme case: $K=1$

24

The leave-one-out estimate

Extreme case: $K=1$

$\mathbb{E}[e_1] = E_{\text{out}}(\mathcal{G}_1)$

$$E_{\text{cv}} = \frac{1}{N} \sum_{n=1}^N e_n$$

Theorem. E_{cv} is an unbiased estimate of $\bar{E}_{\text{out}}(N-1)$.

Expected E_{out} when learning with $N-1$ points.

25

Cross validation

The leave-one-out estimate is expensive to compute!

Cross validation:

- Randomly partition the data into k parts ("folds").
- Set one fold aside for evaluation and train a model on the remaining $k-1$ folds and evaluate it on the held-out fold.
- Repeat until each fold has been used for evaluation

26

Cross validation

The leave-one-out estimate is expensive to compute!

Cross validation:

- Randomly partition the data into k parts ("folds").
- Set one fold aside for evaluation and train a model on the remaining $k-1$ folds and evaluate it on the held-out fold.
- Repeat until each fold has been used for evaluation

- The reported error is the average over the errors for each fold.

27

Cross validation

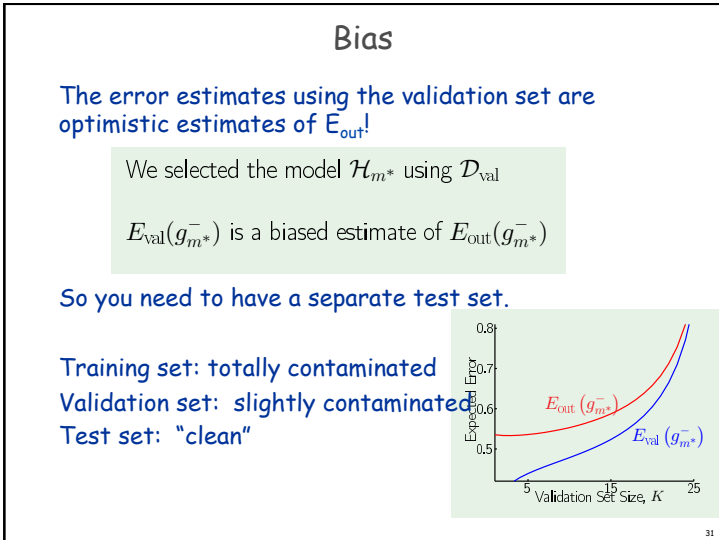
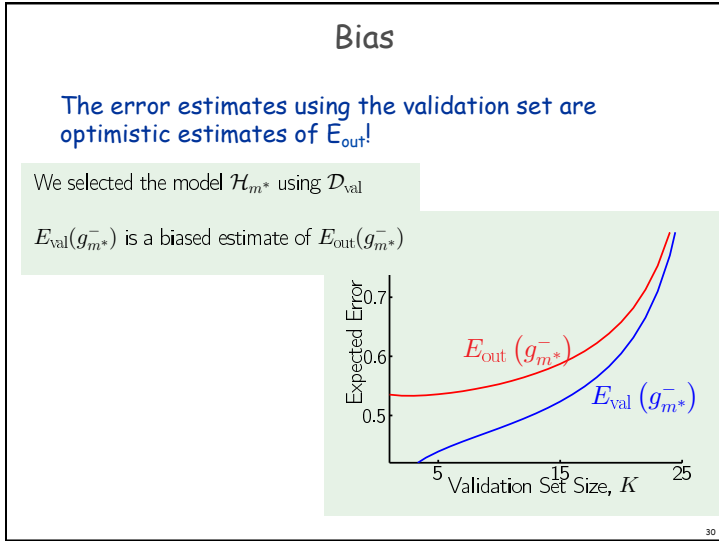
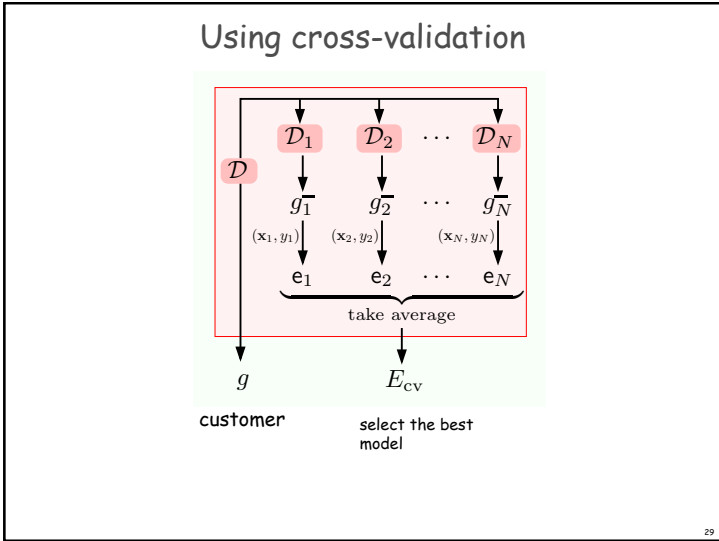
The leave-one-out estimate is expensive to compute!

Cross validation:

- Randomly partition the data into k parts ("folds").
- Set one fold aside for evaluation and train a model on the remaining $k-1$ folds and evaluate it on the held-out fold.
- Repeat until each fold has been used for evaluation

Stratified-cross validation aims at achieving roughly the same class distribution in each fold.

28



Measures of classifier performance

Classifier performance can be summarized by a table known as the **confusion matrix** or contingency table:

		predicted labels:	
		-1	1
true labels	-1	1439	61
	1	62	1438

Measures of classifier performance

Let's take a closer look at the contingency table:

true labels	predicted labels:	
	-1	1
	-1 1439 61	1 62 1438

How do we compute error from the contingency table?

33

Measures of classifier performance

For binary classification problems it is customary to express the contingency table as:

true labels	predicted labels:	
	-1	1
	-1 TN FP	1 FN TP

- TP - number of true positives
- TN - number of true negatives
- FP - number of false positives
- FN - number of false negatives

34

Measures of classifier performance

For binary classification problems it is customary to express the contingency table as:

true labels	predicted labels:		
	-1	1	
	-1 TN FP	1 FN TP	Neg = TN+FP Pos = TP+FN

- True positive rate/sensitivity/recall: TP / Pos
- True negative rate/specificity: TN / Neg
- False positive rate: FP / Neg
- Precision: $TP / (TP + FP)$

35

Measures of classifier performance

Suppose you have a dataset with very few positive examples compared to negative examples (unbalanced data)

A classifier that classifies every example as negative would still attain high accuracy (this is called the majority class classifier).

Need an alternative measure of accuracy!

36

The choice of classification threshold

All the classifiers we will study provide a scoring function whose magnitude indicates how sure we are it belongs to a given class. For example: $w^T x + b$

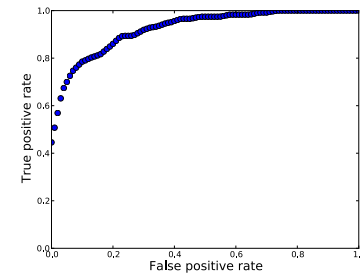
The choice of the threshold is somewhat arbitrary, and in a given application we may prefer to ignore positive predictions that are associated with small scores

To have a view of classifier performance that is independent of the choice of threshold we consider the ROC curve.

37

ROC curve

The ROC curve is a plot of the true positive rate as a function of false positive rate as you vary the classification threshold



How does the ROC curve of a perfect classifier look like?
For a random classifier?

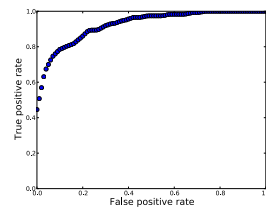
ROC curve computed on the heart disease dataset from the UCI repository

38

ROC curves and ranking

An ROC curve is often summarized by the area under the curve (AUC).

AUC = 0.92



AUC is essentially the probability that a positive example will get a higher score than a negative example

39

ROC curves

This is also a nice way of comparing classifiers:

