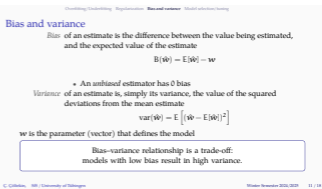
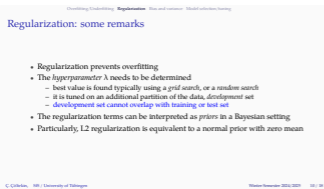
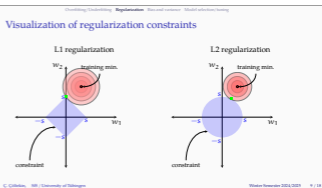
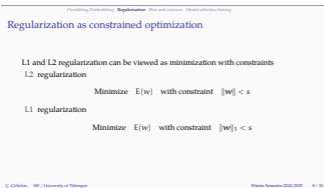
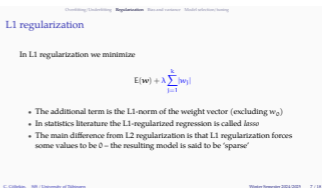
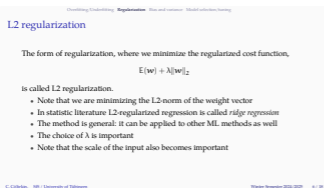
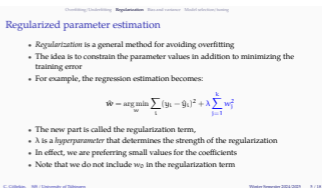
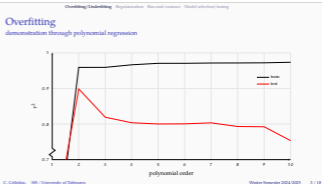
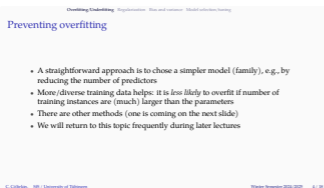


Overfitting & Underfitting

We want our models to *generalize*, perform well on unseen data.

- *Overfitting* occurs when the model learns the idiosyncrasies of the training data
- *Underfitting* occurs when the model is not flexible enough for solving the problem at hand

We want simpler models, but not too simple for the task at hand.



Bias-variance, underfitting-overfitting

- Bias and variance are properties of estimators
- We want estimators with low bias, low variance
- Complex models tend to overfit – and exhibit high variance
- Simple models tend to have low variance, but likely to have (high) bias

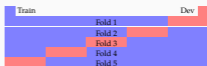
Model selection & hyperparameter tuning

- Our aim is to reduce the error on unseen data
- The evaluation practice should reflect that
- We can estimate the test error on a *development set* (*validation or held-out data*):
 - Split the data at hand as *training and development set*
 - Train alternative models (different hyperparameters) on the training set
 - Choose the model with best development set performance

Cross validation

- To avoid overfitting, we want to tune our models on a *development set*
- But (labeled) data is valuable
- Cross validation is a technique that uses all the data, for both training and tuning with some additional effort
- Besides tuning hyper-parameters, we may also want to get ‘average’ parameter estimates over multiple folds

K-fold Cross validation



- At each fold, we hold part of the data for testing, train the model with the remaining data
- The special case where k equal to the number of data points is called *leave-one-out cross validation*

The choice of k in k -fold CV

- Increasing k
 - reduces the bias: the estimates converge to true value of the measure (e.g., R^2) in the limit
 - increases the variance: smaller held-out sets produce more varied parameter estimates
 - is generally computationally expensive
- 5- or 10-fold cross validation is common practice (and found to have a good balance between bias and variance)

Comparing with a baseline

- The performance measures are only meaningful if we have something to compare against
 - random: does the model do anything useful at all?
 - majority class: does the classifier work better than predicting the majority class all the time?
 - state-of-the-art: how does your model compare against known (non-trivial) models?
- In comparing different models we use another split of the data, *test set*
- Ideally test set is used only once – we want to avoid tuning the system on the test data
- Differences between models are exactly repeatable when the same test set is used (by different studies)
- Differences are reliable if your test set size is large enough
- Use statistical tests when comparing different models/methods

Summary

The first principle is that you must not fool yourself and you are the easiest person to fool. – Richard P. Feynman

- We want models with low bias and low variance
- Evaluating ML system requires special care:
 - Tuning your system on a development set
 - Cross-validation allows efficient use of labeled data during tuning
 - A test set is often used when comparing results obtained by different models

Next:

- Classification

