# Common Mistakes in Machine Learning

## (Not exhaustive)

# Bad Annotation of Training/ Testing Data Set

# Poor Understanding Algorithm Assumptions

# Poor Understanding of algorithm parameters

(use of defaults)

# What is our objective?

# Not understanding the data
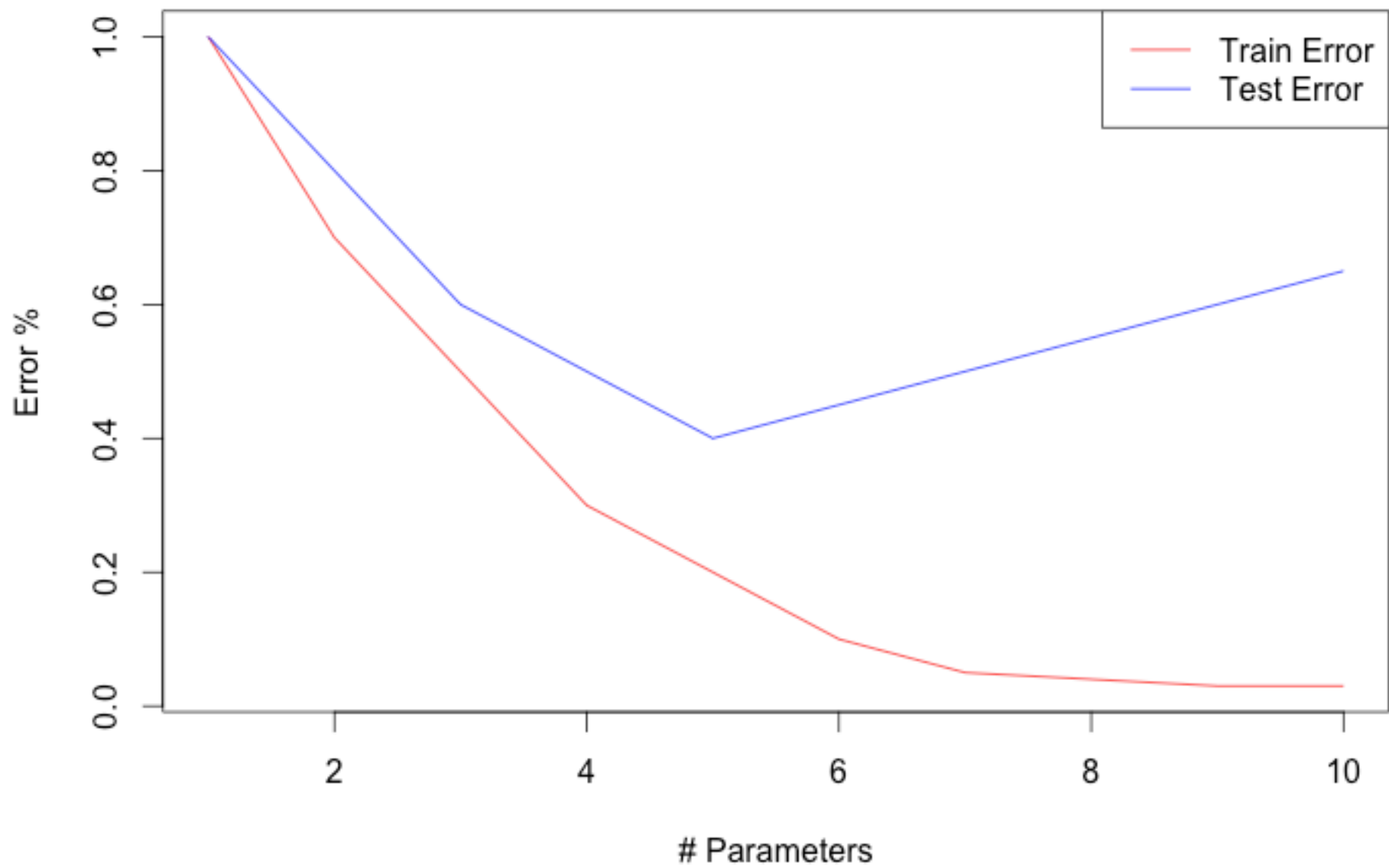
No  EDA - Going into the data blind
(outliers, multicollinearity)
Not understanding population this represents
Insufficient data

# Ignoring n<<P

Potential to overfit

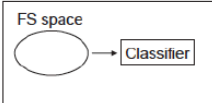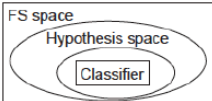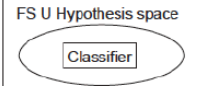# Train/Test Performance

# Allowing Leakage

Features, Information (Classes)

# Motivation for Feature Selection

- Reducing dimensionality
- Improving learning efficiency
- Increasing predicative accuracy
- Reducing complexity of learned results

# Feature Selection Approaches

- **Filter Methods**

- **Wrapper Methods**

- **Embedded Methods**

| | Model search | | Advantages | Disadvantages | Examples |
|---|---|---|---|---|---|
| **Filter** | FS space → Classifier | Univariate | Fast<br>Scalable<br>Independent of the classifier | Ignores feature dependencies<br><br>Ignores interaction with the classifier | Chi-square<br>Euclidean distance<br>t-test<br>Information gain, Gain ratio [6] |
| | | Multivariate | Models feature dependencies<br>Independent of the classifier<br>Better computational complexity<br>than wrapper methods | Slower than univariate techniques<br>Less scalable than univariate<br>techniques<br>Ignores interaction with the classifier | Correlation based feature selection (CFS) [45]<br>Markov blanket filter (MBF) [62]<br>Fast correlation based<br>feature selection (FCBF) [136] |
| **Wrapper** | FS space<br>Hypothesis space<br>Classifier | Deterministic | Simple<br>Interacts with the classifier<br>Models feature dependencies<br>Less computationally intensive<br>than randomized methods | Risk of over fitting<br>More prone than randomized algorithms<br>to getting stuck in a local optimum<br>(greedy search)<br>Classifier dependent selection | Sequential forward selection (SFS) [60]<br>Sequential backward elimination (SBE) [60]<br>Plus $q$ take-away $r$ [33]<br>Beam search [106] |
| | | Randomized | Less prone to local optima<br>Interacts with the classifier<br>Models feature dependencies | Computationally intensive<br>Classifier dependent selection<br>Higher risk of overfitting<br>than deterministic algorithms | Simulated annealing<br>Randomized hill climbing [110]<br>Genetic algorithms [50]<br>Estimation of distribution algorithms [52] |
| **Embedded** | FS ∪ Hypothesis space<br>Classifier | | Interacts with the classifier<br>Better computational complexity<br>than wrapper methods<br>Models feature dependencies | Classifier dependent selection | Decision trees<br>Weighted naive Bayes [28]<br>Feature selection using<br>the weight vector of SVM [44, 125] |

Saeys et al 2005

# Not matching approach to question (or not having a question)

If you try hard enough, you can beat the data into submission to say anything – avoid this **please**.

# Poor Understanding of algorithm
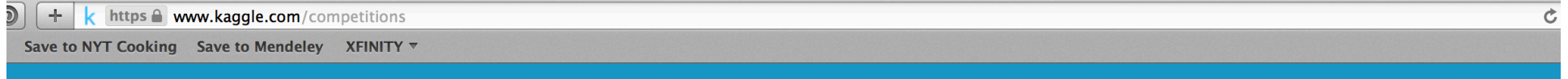
Methods Overview + In-class Discussion

# LINEAR DISCRIMINANT ANALYSIS

Methods Overview + In-class Discussion

# CLASSIFICATION + REGRESSION TREES

# https://www.kaggle.com/

https ⟶ www.kaggle.com/competitions

Save to NYT Cooking    Save to Mendeley    XFINITY ▾

## Active Competitions

**All Competitions**

### Active Competitions

**Ultrasound Nerve Segmentation**
Identify nerve structures in ultrasound images of the neck

**2 months**
**96** teams
**139** scripts
**$100,000**

**Draper Satellite Image Chronology**
Can you put order to space and time?

**31 days**
**261** teams
**341** scripts
**$75,000**

**State Farm Distracted Driver Detection**
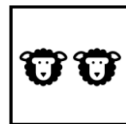Can computer vision spot distracted drivers?

**2 months**
**864** teams
**571** scripts
**$65,000**

**Expedia Hotel Recommendations**
Which hotel type will an Expedia customer book?

**14 days**
**1586** teams
**1890** scripts
**$25,000**

**Avito Duplicate Ads Detection**
Can you detect duplicitous duplicate ads?

**45 days**
**242** teams
**157** scripts
**$20,000**

**40 days**