# HACKERMATH FOR ML

Intro to Stats & Maths for Machine Learning

Amit Kapoor
@amitkaps

Bargava Subramanian
@bargava

# "WHAT I CANNOT CREATE, I DO NOT UNDERSTAND"

Richard Feynman

# PHILOSOPHY OF HACKERMATH

"Hacker literally means developing mastery over something."
Paul Graham

Here we will aim to learn Math essential for Data Science in this hacker way.

# THREE KEY QUESTIONS

» Why do you need to understand the math?

» What math knowledge do you need?

» Why approach it the hacker's way?

# APPROACH

» Understand the Math.

» Code it to learn it.

» Play with code.

# MODULE 1: LINEAR ALGEBRA
## SUPERVISED ML - REGRESSION, CLASSIFICATION

» Solve $Ax = b$ for $n \times n$

» Solve $Ax = b$ for $n \times p + 1$

» Linear Regression

» Ridge Regularization (L2)

» Bootstrapping

» Logistic Regression (Classification)

# MODULE 2: STATISTICS
## HYPOTHESIS TESTING: A/B TESTING

» Basic Statistics

» Distributions

» Shuffling

» Bootstrapping & Simulation

» A/B Testing

# MODULE 3: LINEAR ALGEBRA CONTD.
## UNSUPERVISED ML: DIMENSIONALITY REDUCTION

» Solve $Ax = \lambda x$ for $n \times n$

» Eigenvectors & Eigenvalues

» Principle Component Analysis

» Cluster Analysis (K-Means)

# SCHEDULE

```
0900 - 1000: Breakfast
1000 - 1130: Session 1
1130 - 1145: Tea Break
1145 - 1315: Session 2
1315 - 1400: Lunch
1400 - 1530: Session 3
1530 - 1545: Tea Break
1545 - 1700: Session 4
```

# "IT'S TOUGH TO MAKE PREDICTIONS, ESPECIALLY ABOUT THE FUTURE."

Yogi Berra

# WHAT IS MACHINE LEARNING (ML)?

"[Machine learning is the] field of study that gives computers the ability to learn without being explicitly programmed."
Arthur Samuel

"Machine learning is the study of computer algorithm that improve automatically through experience"
Tom Mitchell

# ML PROBLEMS

» "Is this cancer?"

» "What is the market value of this house?"

» "Which of these people are friends?"

» "Will this person like this movie?"

» "Who is this?"

» "What did you say?"

» "How do you fly this thing?".

# ML IN USE EVERYDAY

» Search

» Photo Tagging

» Spam Filtering

» Recommendation

» ...

# BROAD ML APPLICATION

» Database Mining e.g. Clickstream data, Business data

» Automating e.g. Handwriting, Natural Language Processing, Computer Vision

» Self Customising Program e.g. Recommendations

# ML THOUGHT PROCESS

# LEARNING PARADIGM

» Supervised Learning

» Unsupervised Learning

» Reinforcement Learning

» Online Learning

# SUPERVISED LEARNING

» Regression

» Classification

# UNSUPERVISED LEARNING

» Clustering

» Dimensionality Reduction

# ML PIPELINE

» Frame: Problem definition

» Acquire: Data ingestion

» Refine: Data wrangline

» Transform: Feature creation

» Explore: Feature selection

» Model: Model creation & assessment

» Insight: Communication

# LINEAR REGRESSION

# LINEAR RELATIONSHIP

$$y_i = \alpha + \beta_1 x_1 + \beta_2 x_2 + ..$$

# OBJECTIVE FUNCTION

$$\epsilon = \sum_{k=1}^{n} (y_i - \hat{y}_i)^2$$

Interactive Example: http://setosa.io/ev/

# LOGIT FUNCTION

$$\sigma(t) = \frac{e^t}{e^t + 1} = \frac{1}{1 + e^{-t}}$$

# LOGISTIC REGRESSION



$$y = b_0 + b_1 x \quad \leftarrow \text{Linear Model}$$

Logistic Model

$$p = \frac{1}{1 + e^{-(b_0 + b_1 x)}}$$

# LOGISTIC RELATIONSHIP

Find the $\beta$ parameters that best fit:
$y = 1$ if $\beta_0 + \beta_1 x + \epsilon > 0$
$y = 0$, otherwise

Follows:

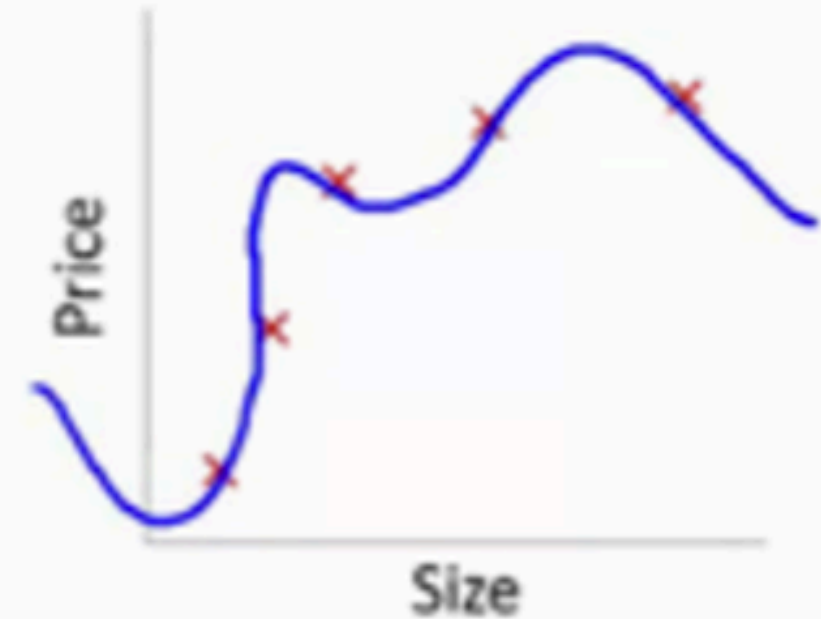$$P(x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

# FITTING A MODEL



High bias
(underfit)

$$\theta_0 + \theta_1 x$$

"Just right"

$$\theta_0 + \theta_1 x + \theta_2 x^2$$

High variance
(overfit)

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$
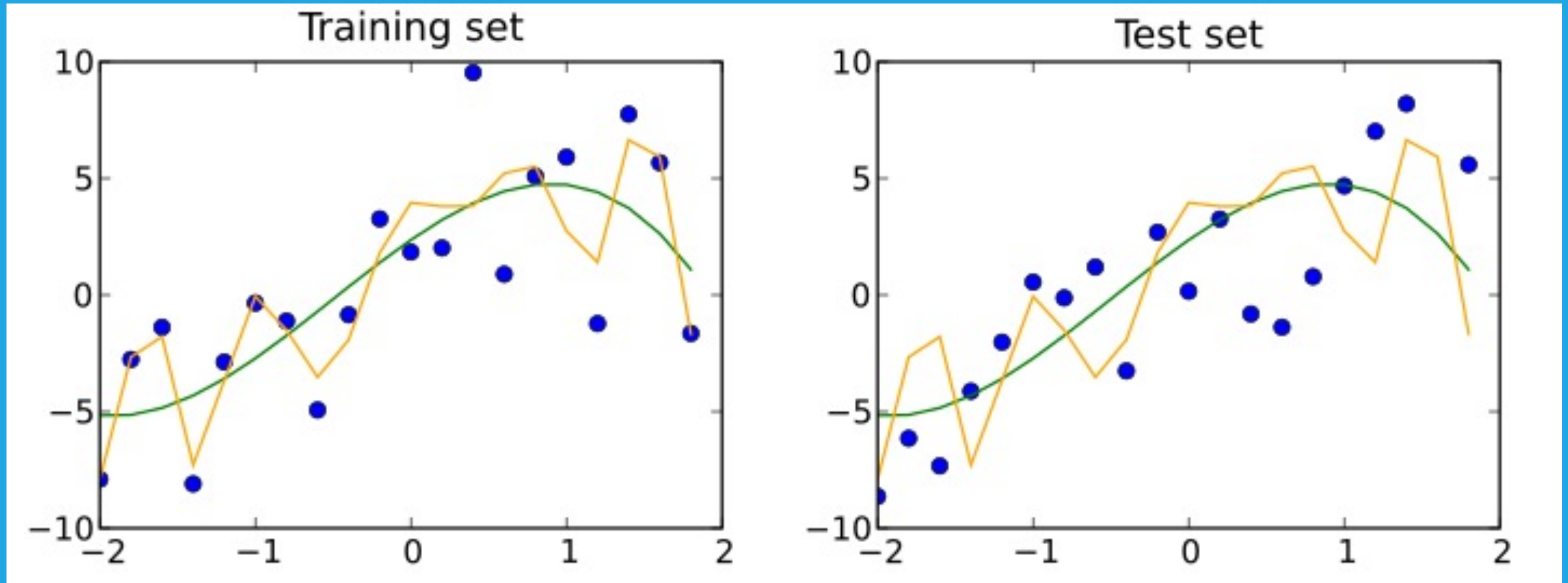
# BIAS-VARIANCE TRADEOFF

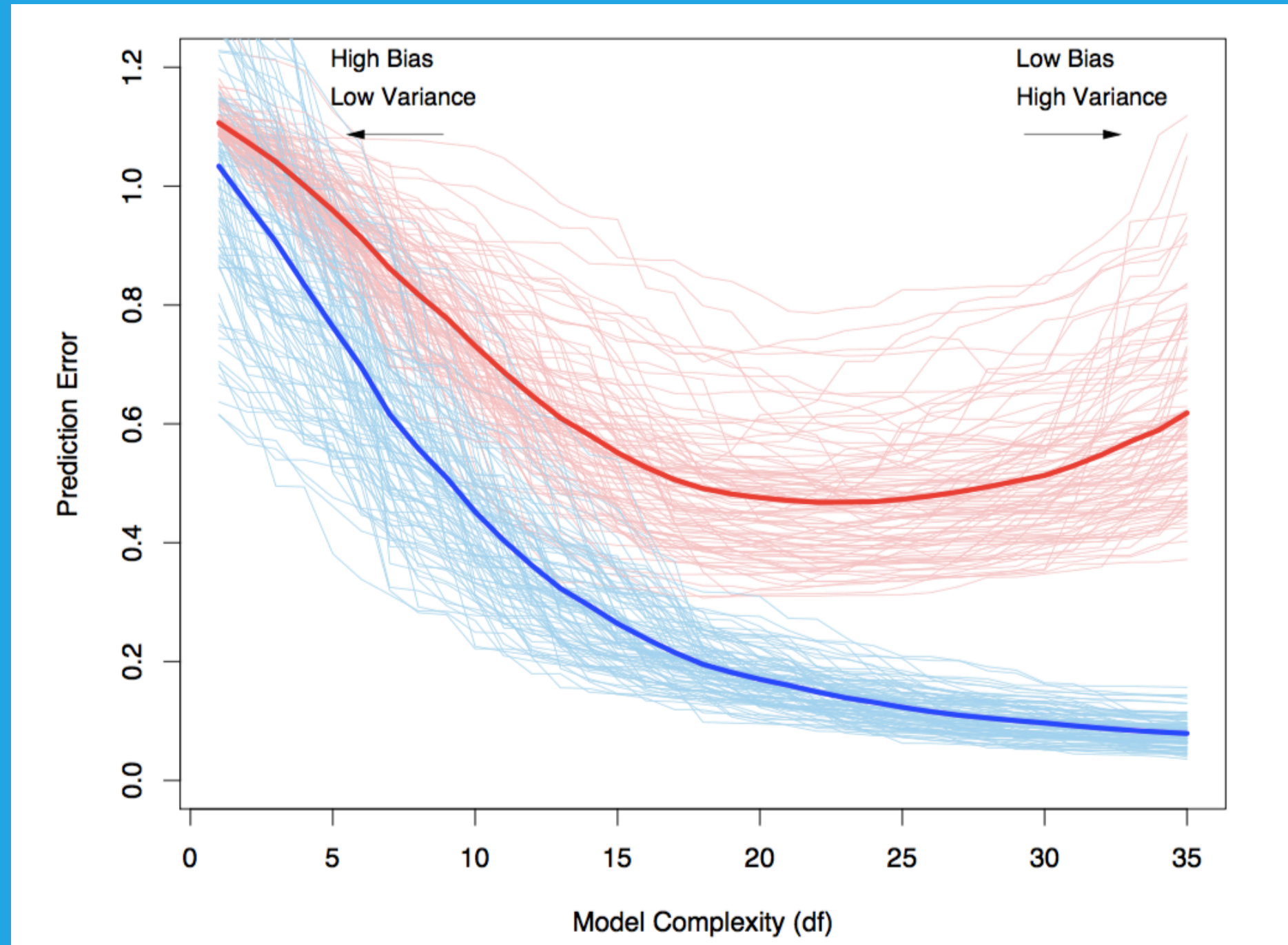# TRAIN AND TEST DATASETS

Split the Data - 80% / 20%

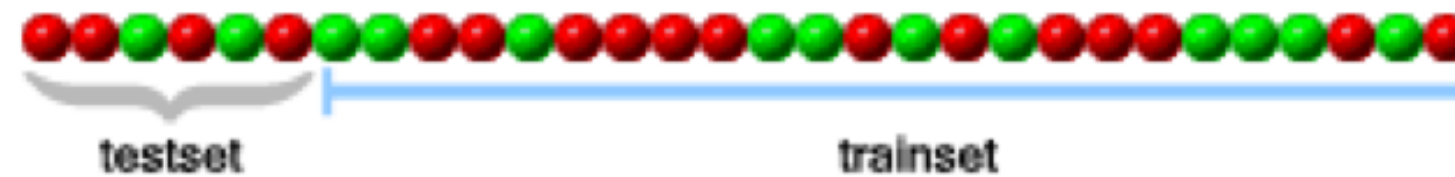# TRAIN AND TEST DATASETS

Measure the error on Test data

# MODEL COMPLEXITY

# CROSS VALIDATION
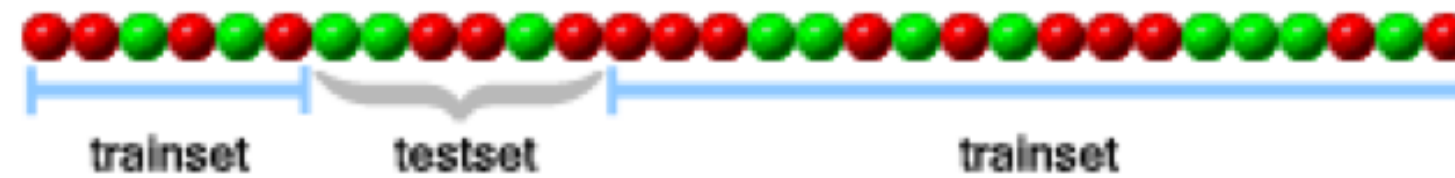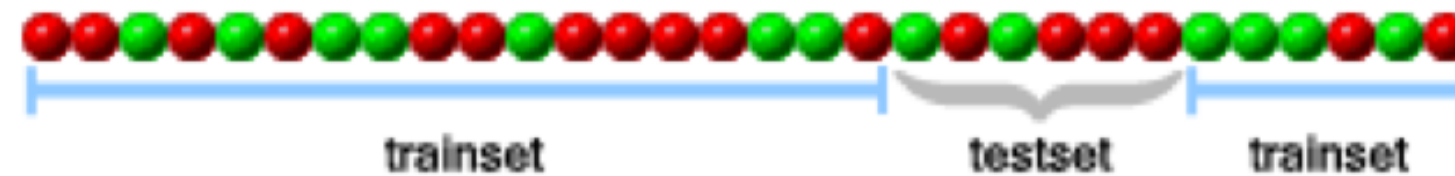


One iteration of a 5-fold Cross-Validation:

1-st fold: testset / trainset

2-nd fold: trainset / testset / trainset

3-rd fold: trainset / testset / trainset
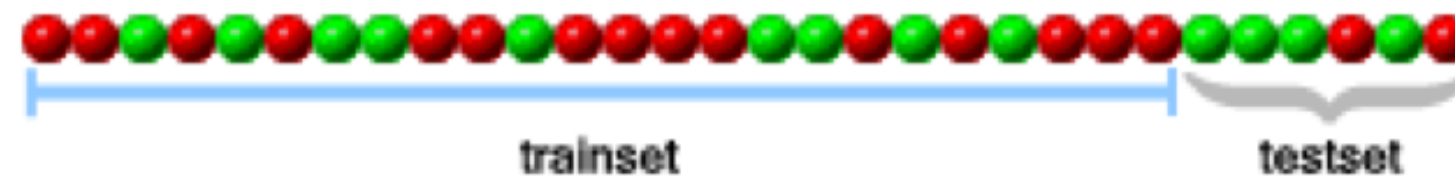
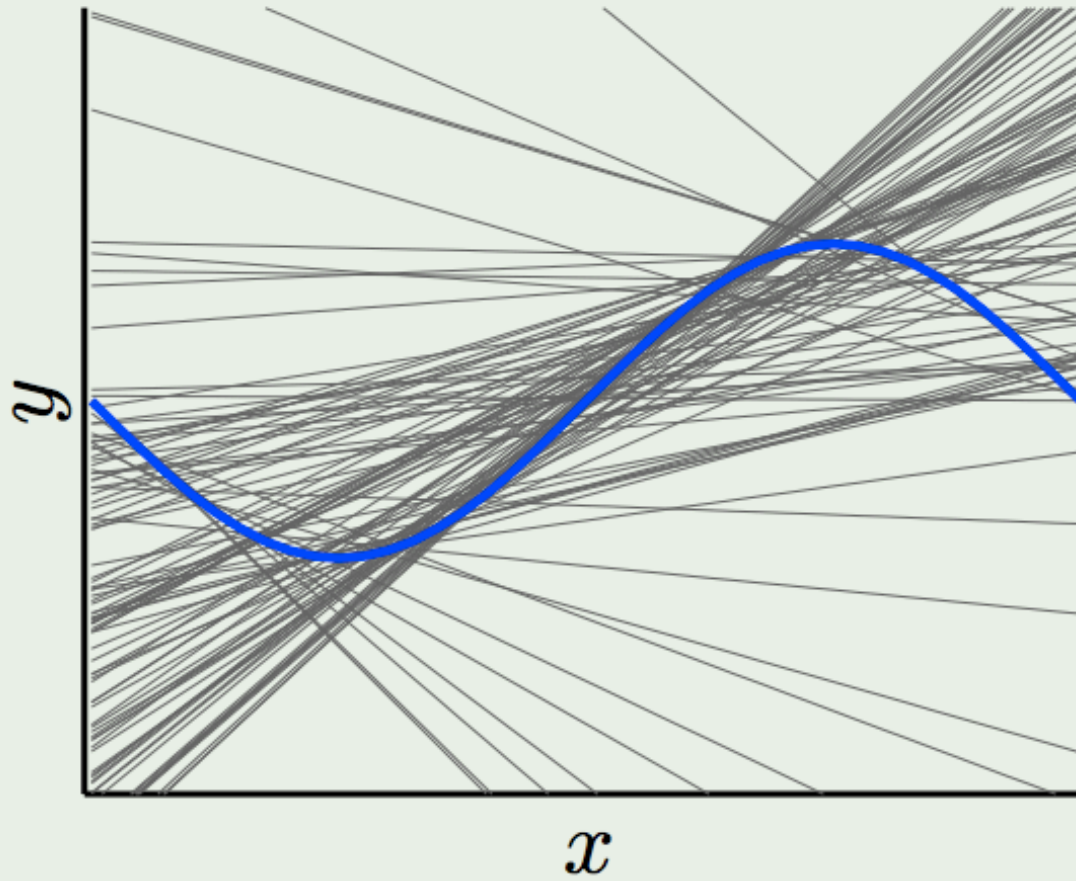4-th fold: trainset / testset / trainset
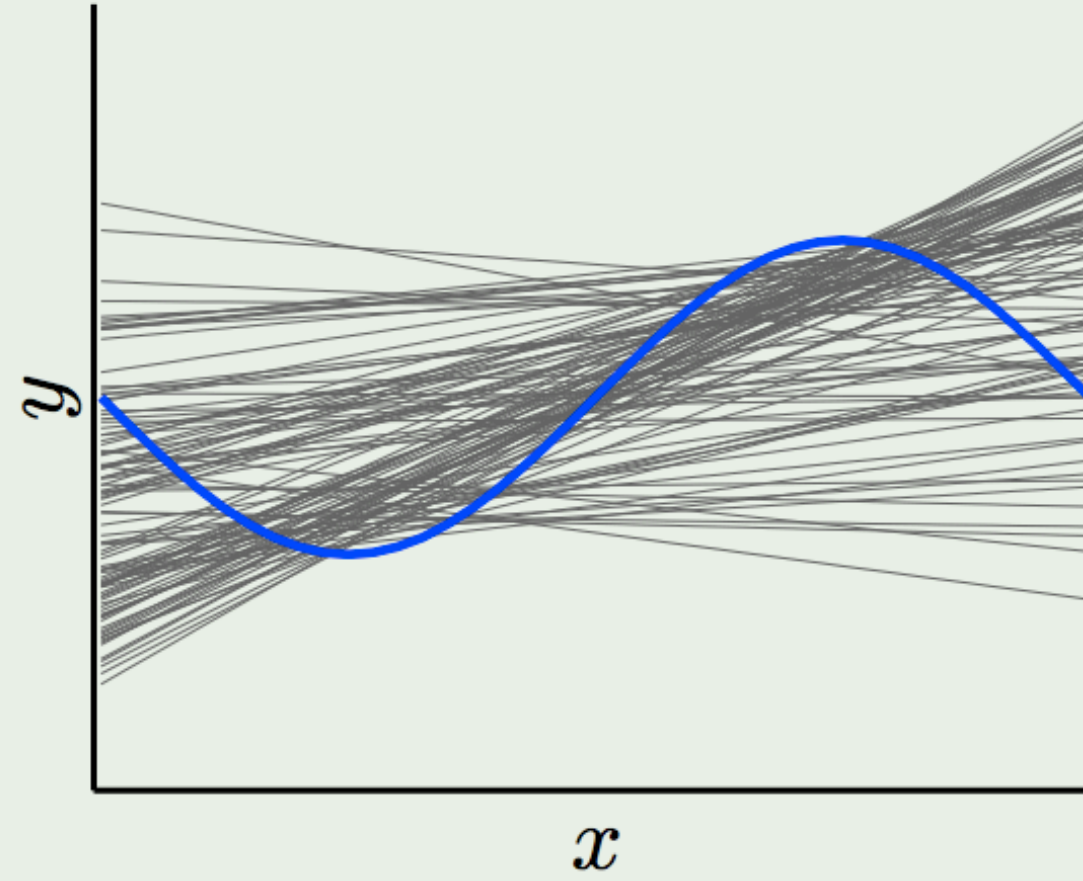
5-th fold: trainset / testset

# REGULARIZATION

Attempts to impose Occam's razor on the solution



without regularization

with regularization

# MODEL EVALUATION

Mean Squared Error

$$MSE = 1/n \sum_{k=1}^{n} (y_i - \hat{y}_i)^2$$

# MODEL EVALUATION

Confusion Matrix

| n=165 | Predicted: NO | Predicted: YES | |
|---|---|---|---|
| Actual: NO | TN = 50 | FP = 10 | 60 |
| Actual: YES | FN = 5 | TP = 100 | 105 |
| | 55 | 110 | |

# MODEL EVALUATION

## Classification Metrics

Recall (TPR) = TP / (TP + FN)

Precision = TP / (TP + FP)

Specificity (TNR) = TN / (TN + FP)



relevant elements

false negatives          true negatives

true positives    false positives

selected elements

How many selected items are relevant?

How many relevant items are selected?

Precision =

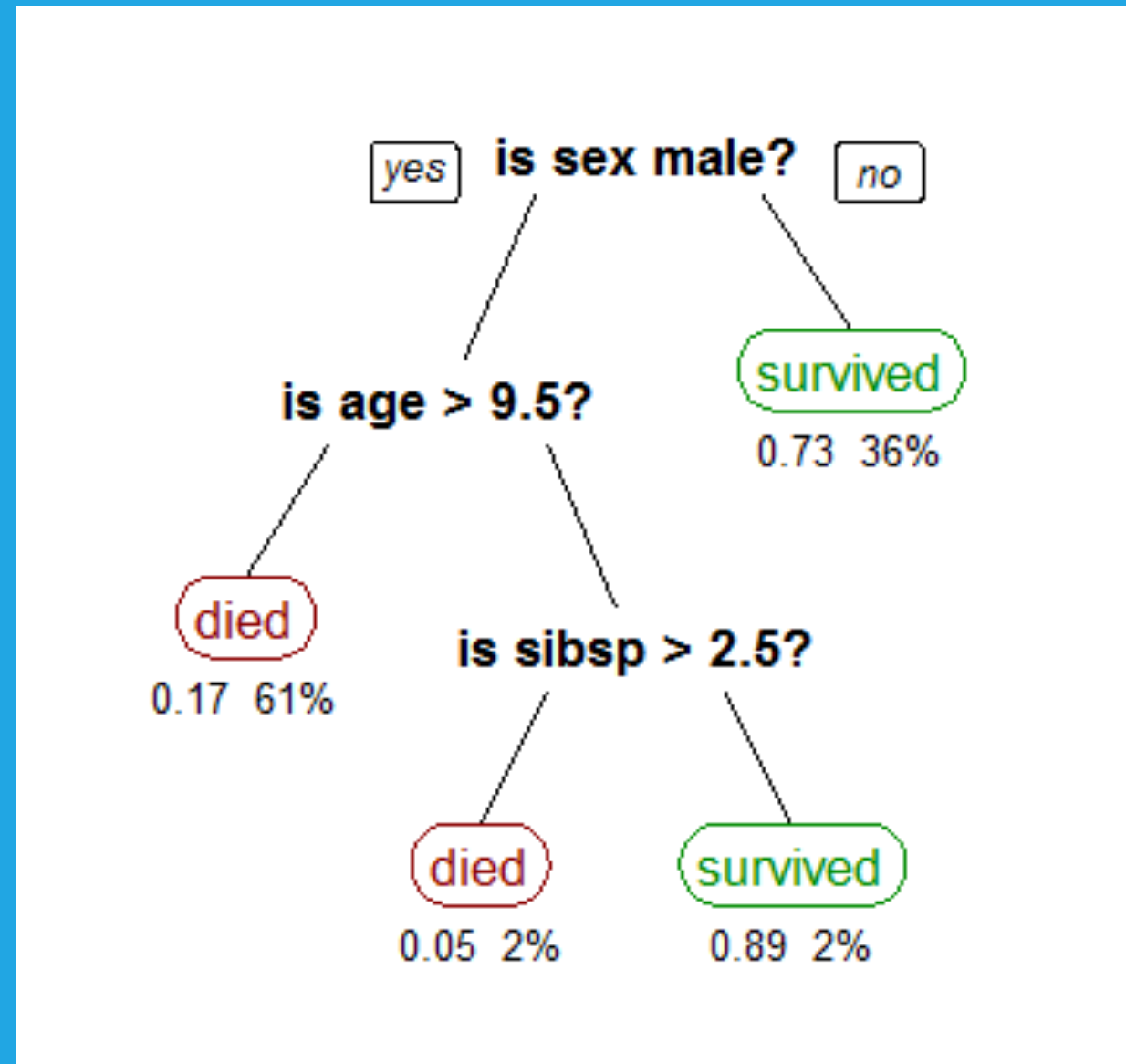Recall =

# MODEL EVALUATION

Receiver Operating
Characteristic Curve

Plot of TPR vs FPR at
different discrimination
threshold

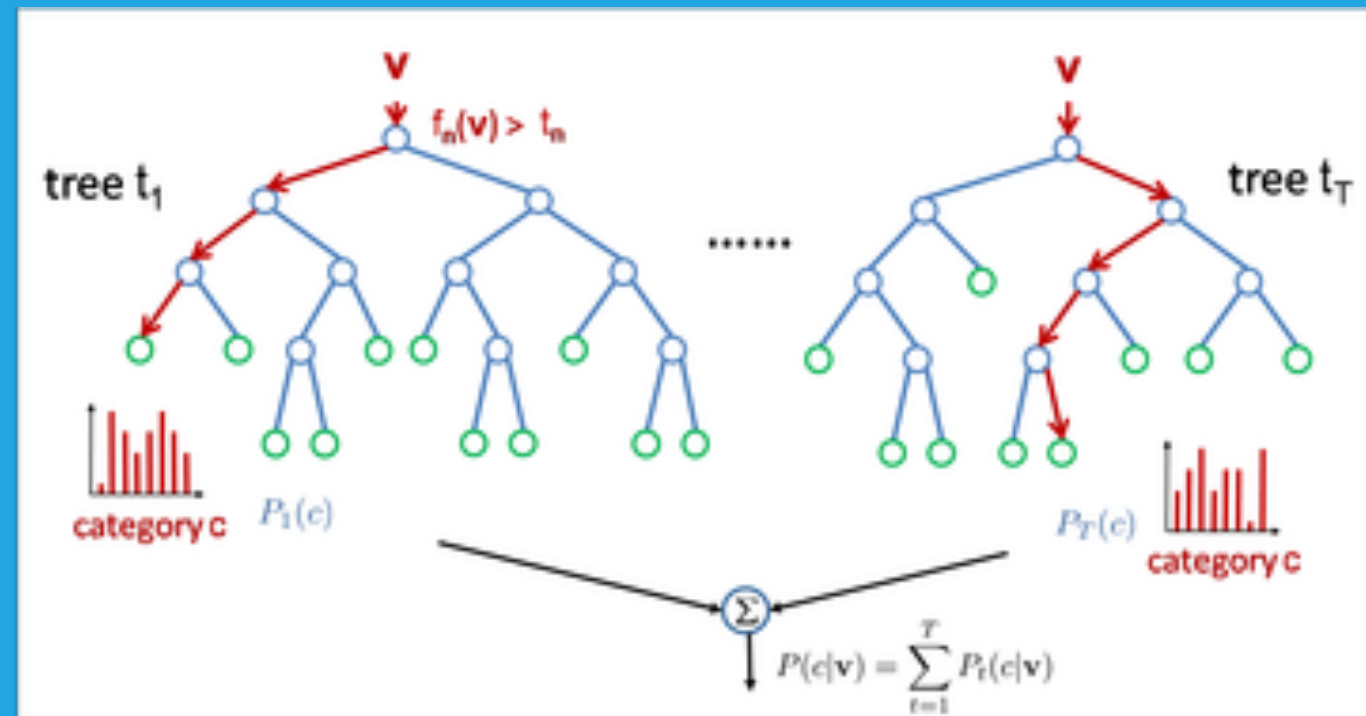# DECISION TREE

Example: Survivor on Titanic

# DECISION TREE

» Easy to interpret

» Little data preparation

» Scales well with data

» White-box model

» Instability — changing variables, altering sequence

» Overfitting

# BAGGING

» Also called bootstrap aggregation, reduces variance

» Uses decision trees and uses a model averaging approach

# RANDOM FOREST

» Combines bagging idea and random selection of features.

» Similar to decision trees are constructed – but at each split, a random subset of features is used.

# "IF YOU TORTURE THE DATA ENOUGH, IT WILL CONFESS."

Ronald Case

# CHALLENGES

» Data Snooping

» Selection Bias

» Survivor Bias

» Omitted Variable Bias

» Black-box model Vs White-Box model

» Adherence to regulations