# INTRODUCTION TO DATA SCIENCE

This lecture is
based on course by E. Fox and C. Guestrin, Univ of Washington
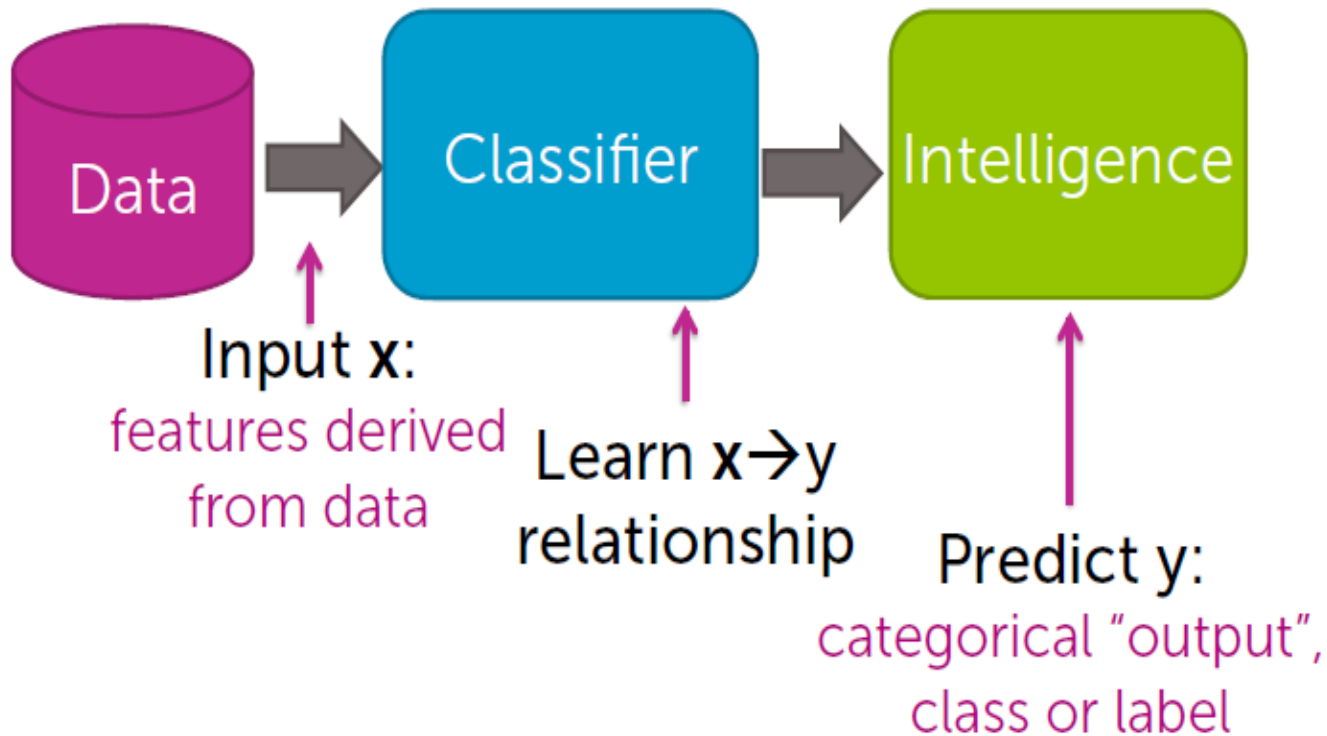
14/11/2017

WFAiS UJ, Informatyka Stosowana
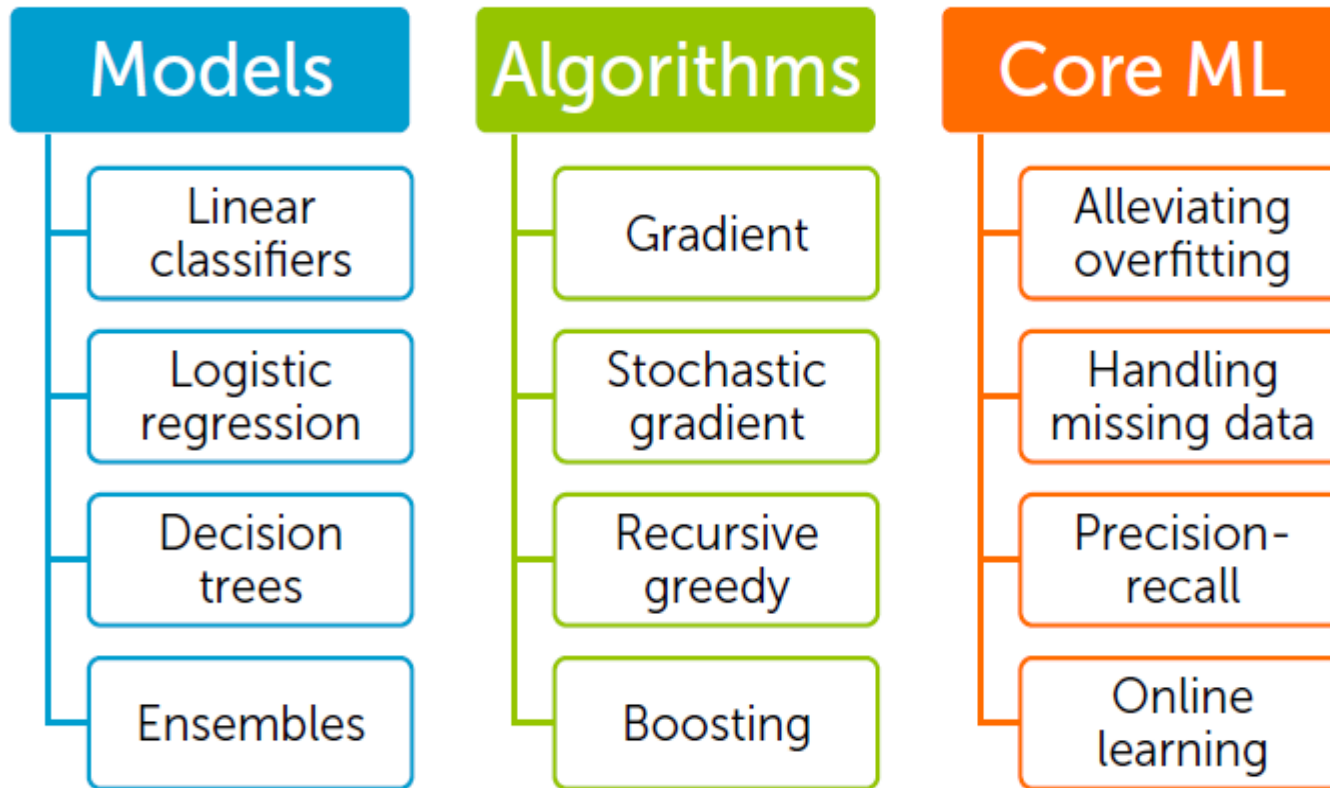II stopień studiów

# What is classification?

From features to predictions

Data → Classifier → Intelligence

Input **x**:
features derived
from data

Learn **x→y**
relationship

Predict y:
categorical "output",
class or label

# Overwiew of content

| Models | Algorithms | Core ML |
|---|---|---|
| Linear classifiers | Gradient | Alleviating overfitting |
| Logistic regression | Stochastic gradient | Handling missing data |
| Decision trees | Recursive greedy | Precision-recall |
| Ensembles | Boosting | Online learning |

14/11/2017

# Sentiment classifier

Input **x:**    Easily best sushi in Seattle.



Sentence Sentiment Classifier

Output: y
Sentiment

14/11/2017

# Classifier

$\hat{y} = +1$

Sentence from review

Input: **x**

Classifier MODEL

Output: y
Predicted class

$\hat{y} = -1$

Note: we'll start talking about 2 classes, and address multiclass later

# Linear classifiers

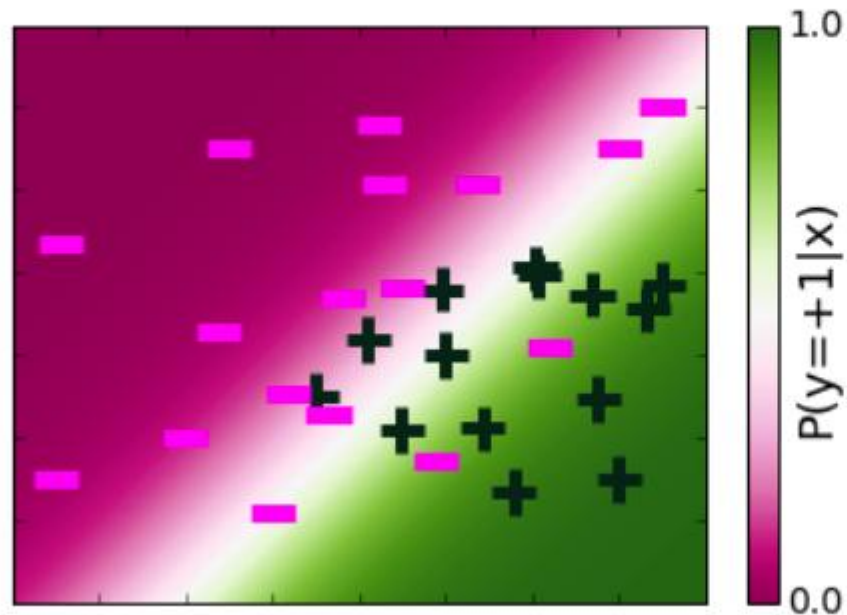| Word | Coefficient |
|------|-------------|
| #awesome | 1.0 |
| #awful | -1.5 |

$\Rightarrow$ Score(x) = 1.0 #awesome − 1.5 #awful
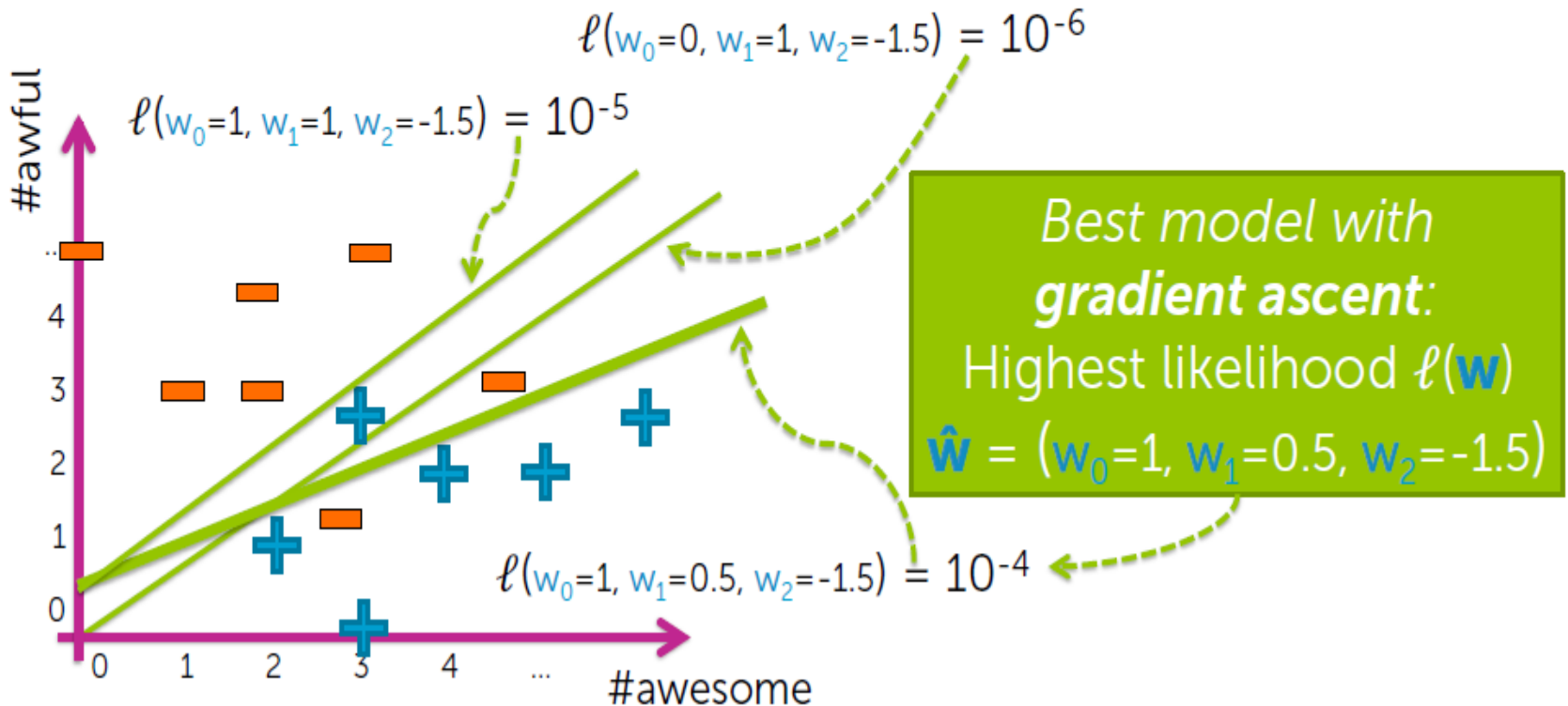
# Logistic regression represents probabilities

$$\hat{P}(y=+1|\mathbf{x},\hat{\mathbf{w}}) = \frac{1}{1 + e^{-\hat{\mathbf{w}}^{T}h(\mathbf{x})}}$$
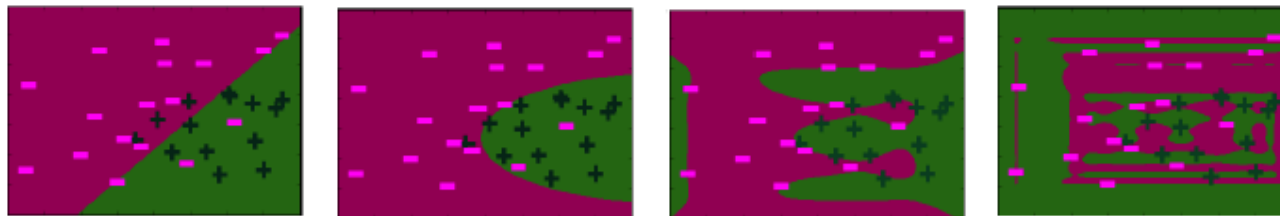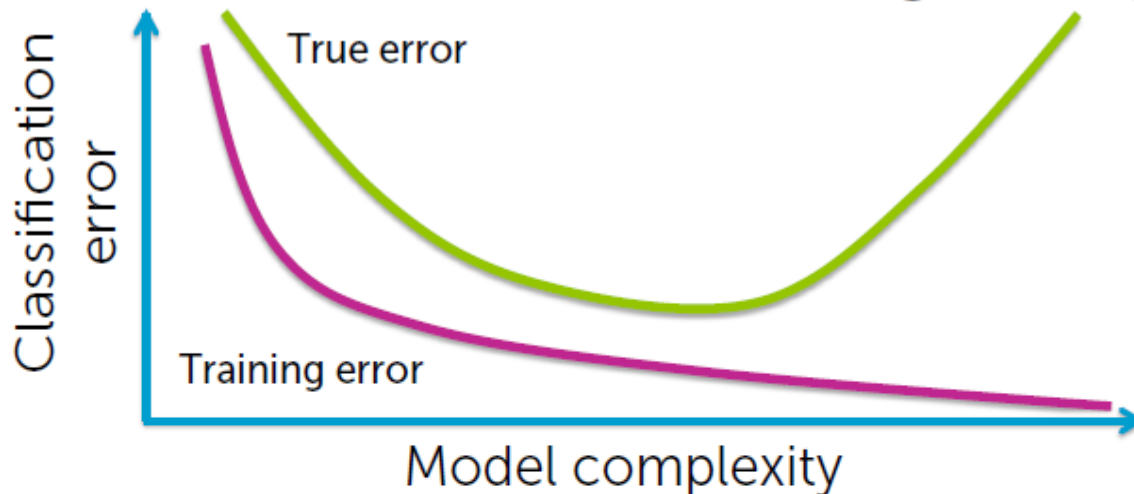


14/11/2017

# Learning „best" classifier

Maximize likelihood over all possible $w_0, w_1, w_2$

$\ell(w_0=0, w_1=1, w_2=-1.5) = 10^{-6}$

$\ell(w_0=1, w_1=1, w_2=-1.5) = 10^{-5}$

#awful

4
3
2
1
0

0   1   2   3   4   ...   #awesome

$\ell(w_0=1, w_1=0.5, w_2=-1.5) = 10^{-4}$

Best model with **gradient ascent**:
Highest likelihood $\ell(\mathbf{w})$
$\hat{\mathbf{w}} = (w_0=1, w_1=0.5, w_2=-1.5)$
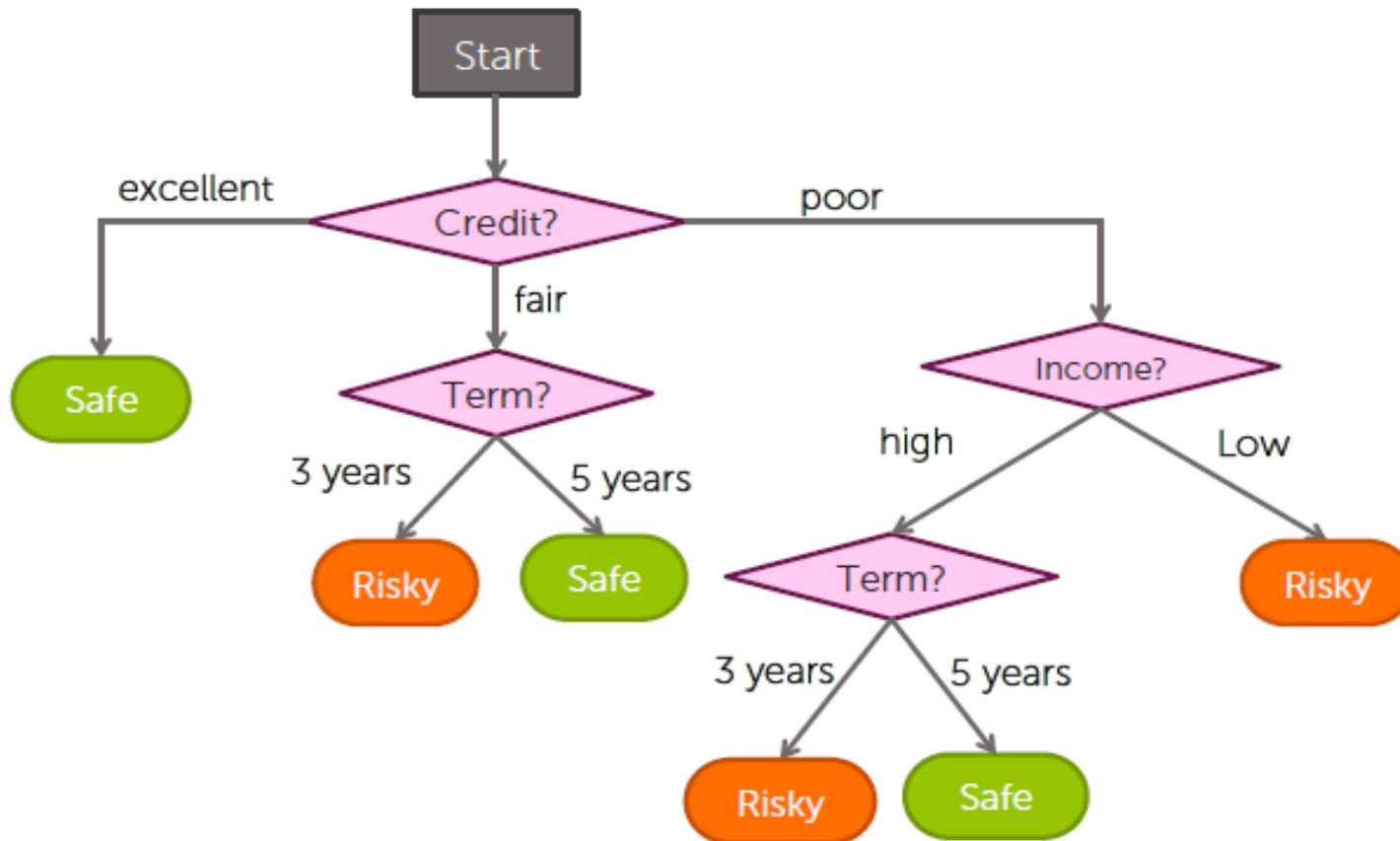
14/11/2017

# Overfitting & regularisation

Use regularization penalty to mitigate overfitting $\quad \ell(\mathbf{w}) - \lambda \|\mathbf{w}\|_2^2$
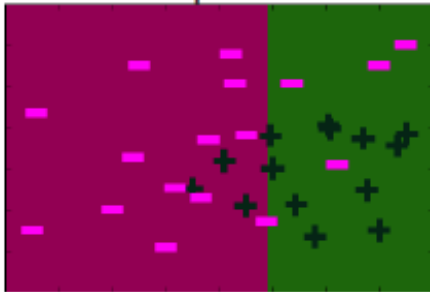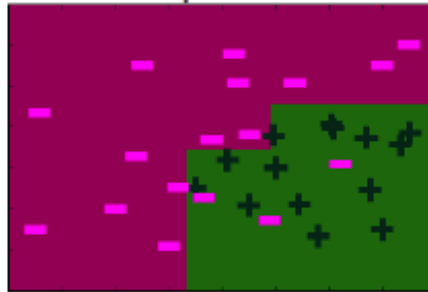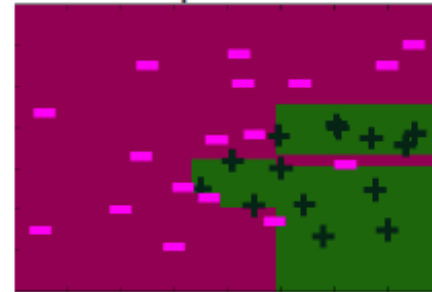
# Decision trees

# Overfitting & decision trees
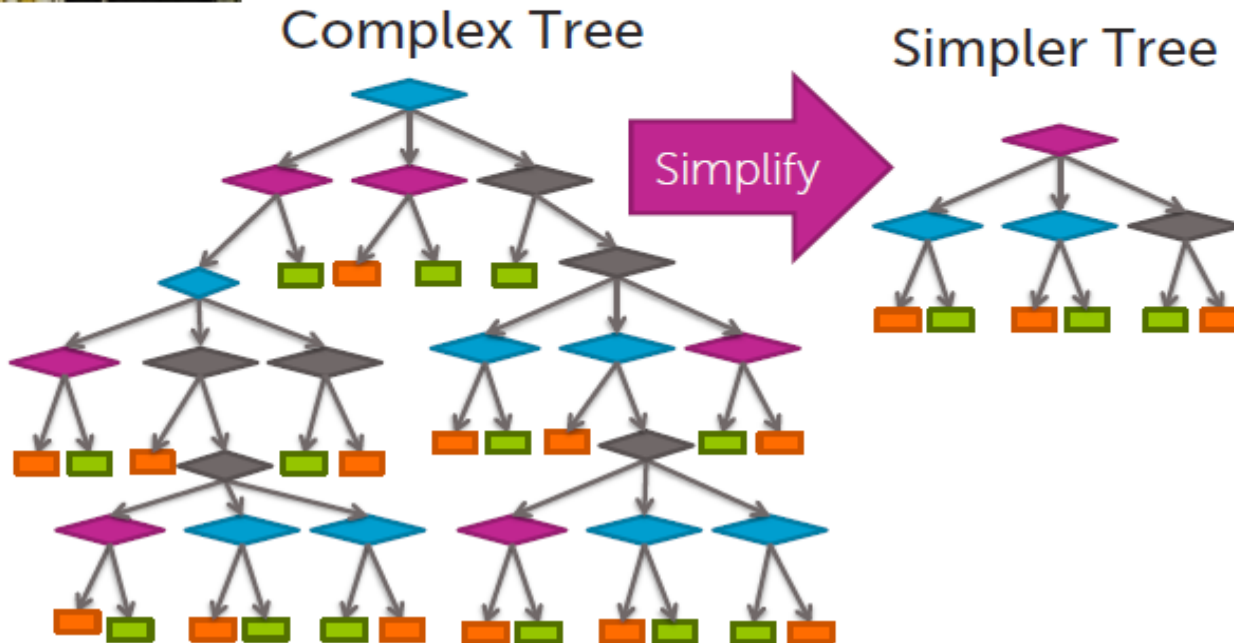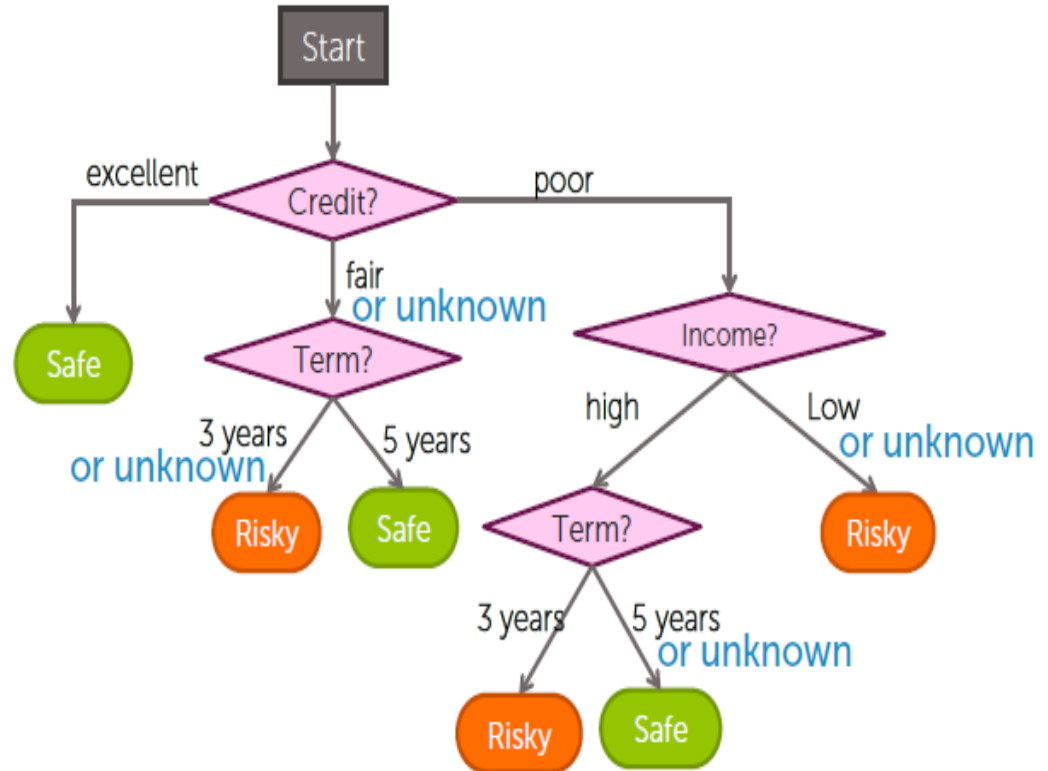
# Alleviate overfitting by learning simpler trees

Occam's Razor: "Among competing hypotheses, the one with fewest assumptions should be selected", William of Occam, 13th Century

**Complex Tree**

Simplify

**Simpler Tree**

# Handling missing data

14/11/2017

# Boosting questions

"Can a set of weak learners be combined to create a stronger learner?" *Kearns and Valiant (1988)*
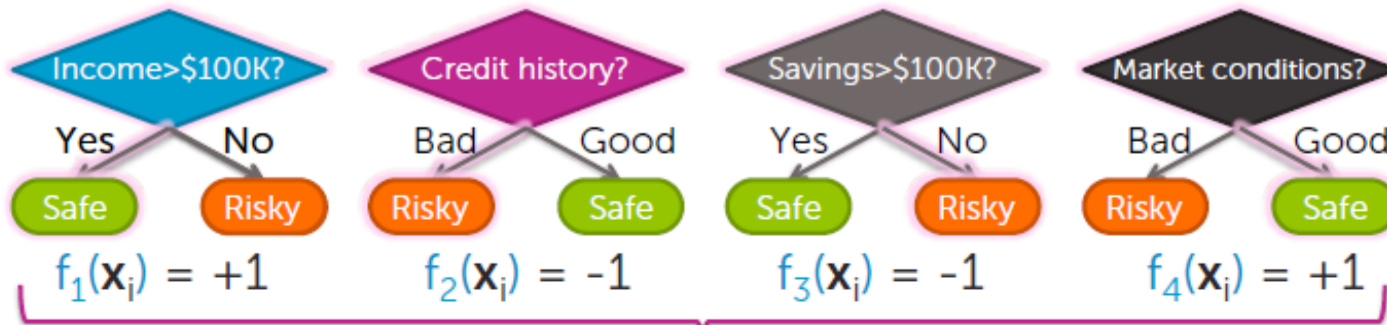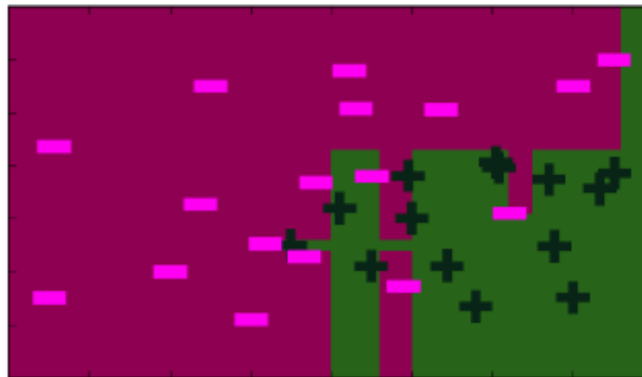
Yes! *Schapire (1990)*

Boosting

Amazing impact: • simple approach • widely used in industry • wins most Kaggle competitions

14/11/2017

# Boosting using AdaBoost

Income>$100K?
Yes → Safe   No → Risky
$f_1(\mathbf{x}_i) = +1$

Credit history?
Bad → Risky   Good → Safe
$f_2(\mathbf{x}_i) = -1$

Savings>$100K?
Yes → Safe   No → Risky
$f_3(\mathbf{x}_i) = -1$

Market conditions?
Bad → Risky   Good → Safe
$f_4(\mathbf{x}_i) = +1$

**Ensemble**: Combine votes from many simple classifiers to learn complex classifiers

# Precision - recall

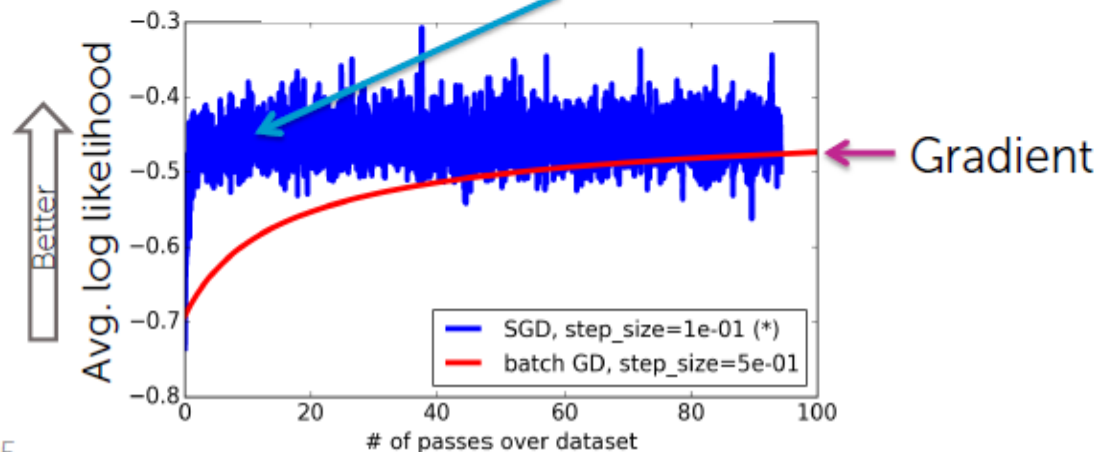# Scalling to huge dataset & on-line learning

4.8B webpages    500M Tweets/day    5B views/day

**Stochastic gradient:** tiny modification to gradient, a lot faster, but annoying in practice



Gradient

Better

Avg. log likelihood

- SGD, step_size=1e-01 (*)
- batch GD, step_size=5e-01

# of passes over dataset

14/11/2017