# DAY 2

## Intelligent Audio Systems:
## A review of the foundations and applications of semantic audio analysis and music information retrieval

*Jay LeBoeuf*
*Imagine Research*
*jay{at}imagine-research.com*

*Rebecca Fiebrink*
*Princeton University*
*rfiebrink{at}princeton.edu*

*July 2010*

These lecture notes contain hyperlinks to the CCRMA Wiki.

On these pages, you can find supplemental material for lectures - providing extra tutorials, support, references for further reading, or demonstration code snippets for those interested in a given topic .

Click on the ⓘ symbol on the lower-left corner of a slide to access additional resources.
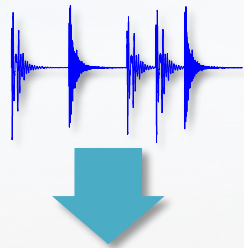
# WIKI REFERENCES...

# Review from Day 1

- What are the 3 major components of a MIR system?
- Name 3 ways of segmenting audio into frames

- What problems did you experience in the lab?
- Follow-up questions?
- Did you try other audio files?
- Did you do the simple instrument recognition?

# FEATURE DEMOS

- Simple re-ordering or slices:
  - Slice up loop into segments and sort via features
  - Play audio
  - Play whole song snippet

# Basic system overview



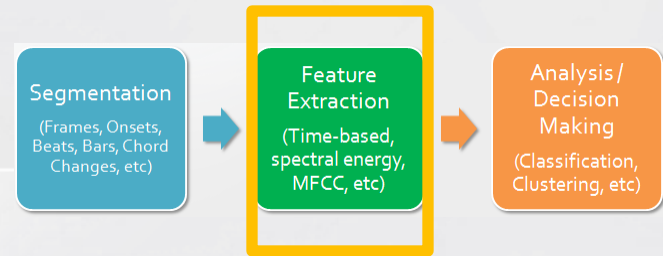Segmentation (Frames, Onsets, Beats, Bars, Chord Changes, etc) → Feature Extraction (Time-based, spectral energy, MFCC, etc) → Analysis / Decision Making (Classification, Clustering, etc)
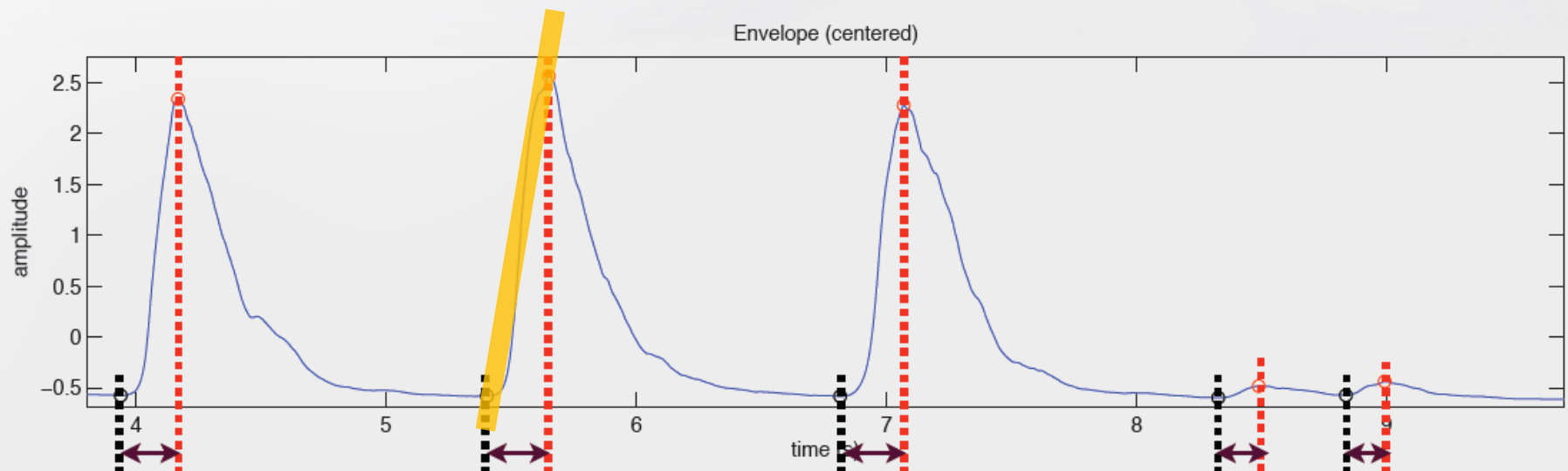
# FEATURE EXTRACTION

# Temporal Information

- Rise time or Attack time- time interval between the onset and instant of maximal amplitude
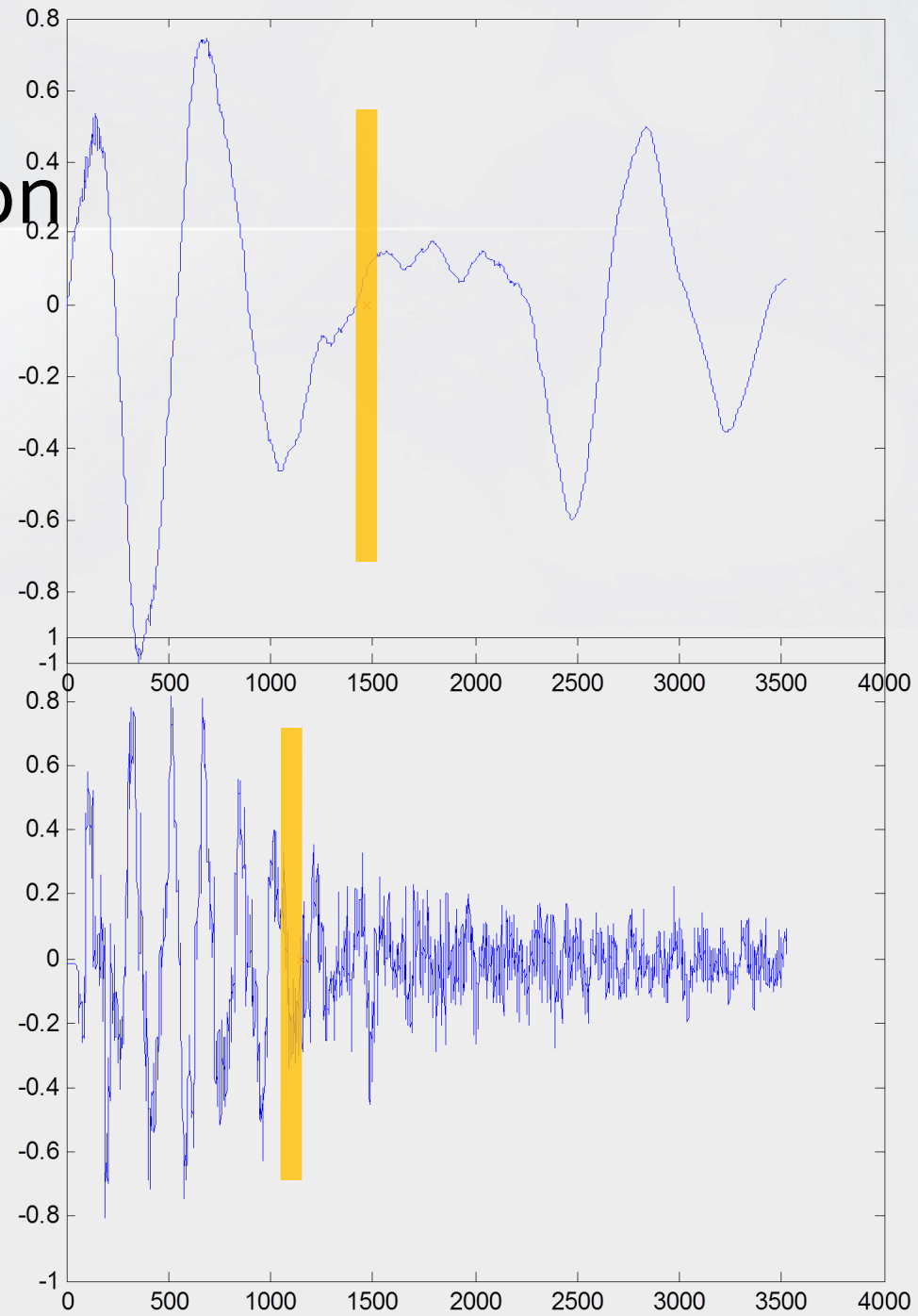- Attack slope



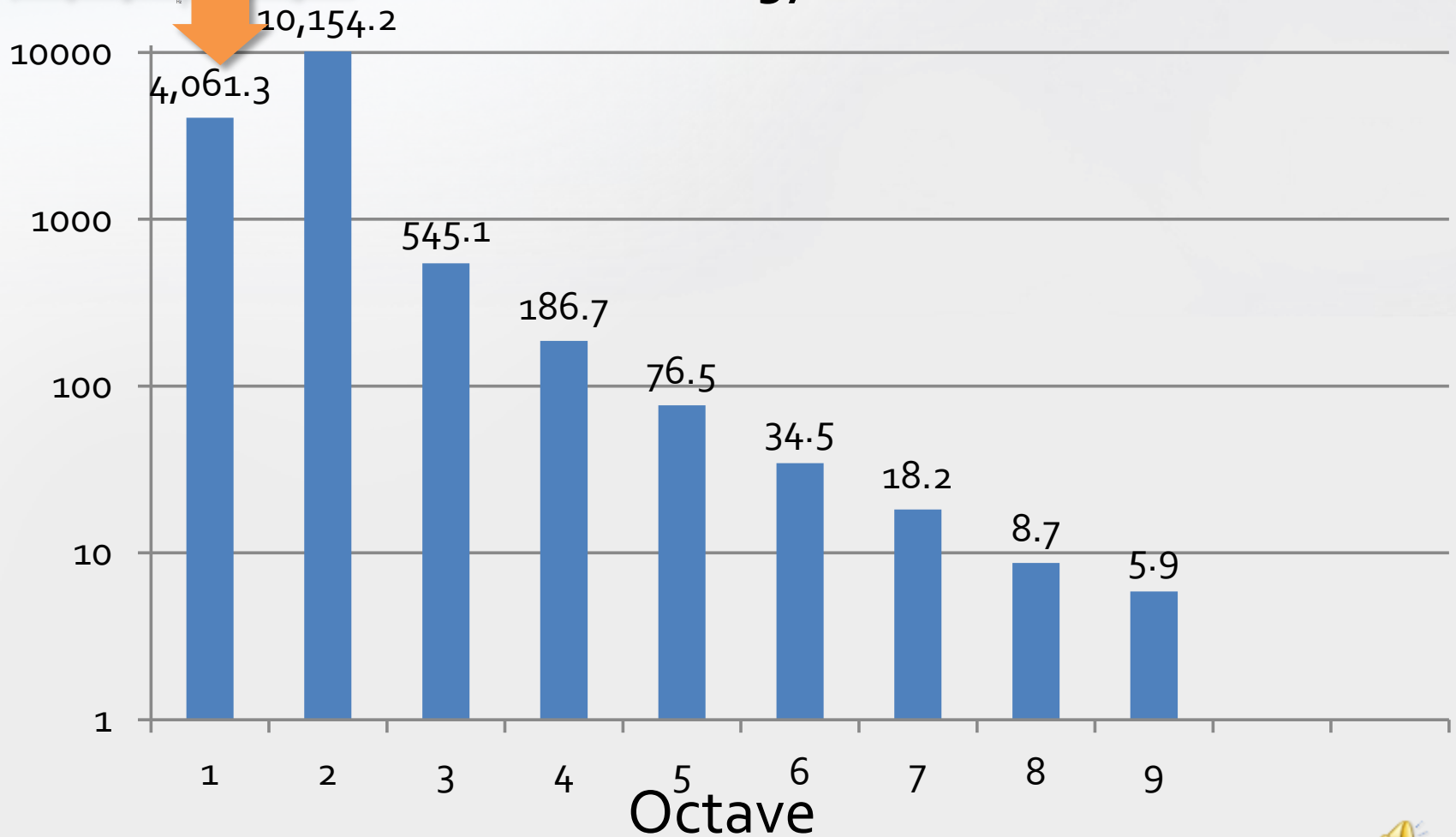Picture courtesy: Olivier Lartillot

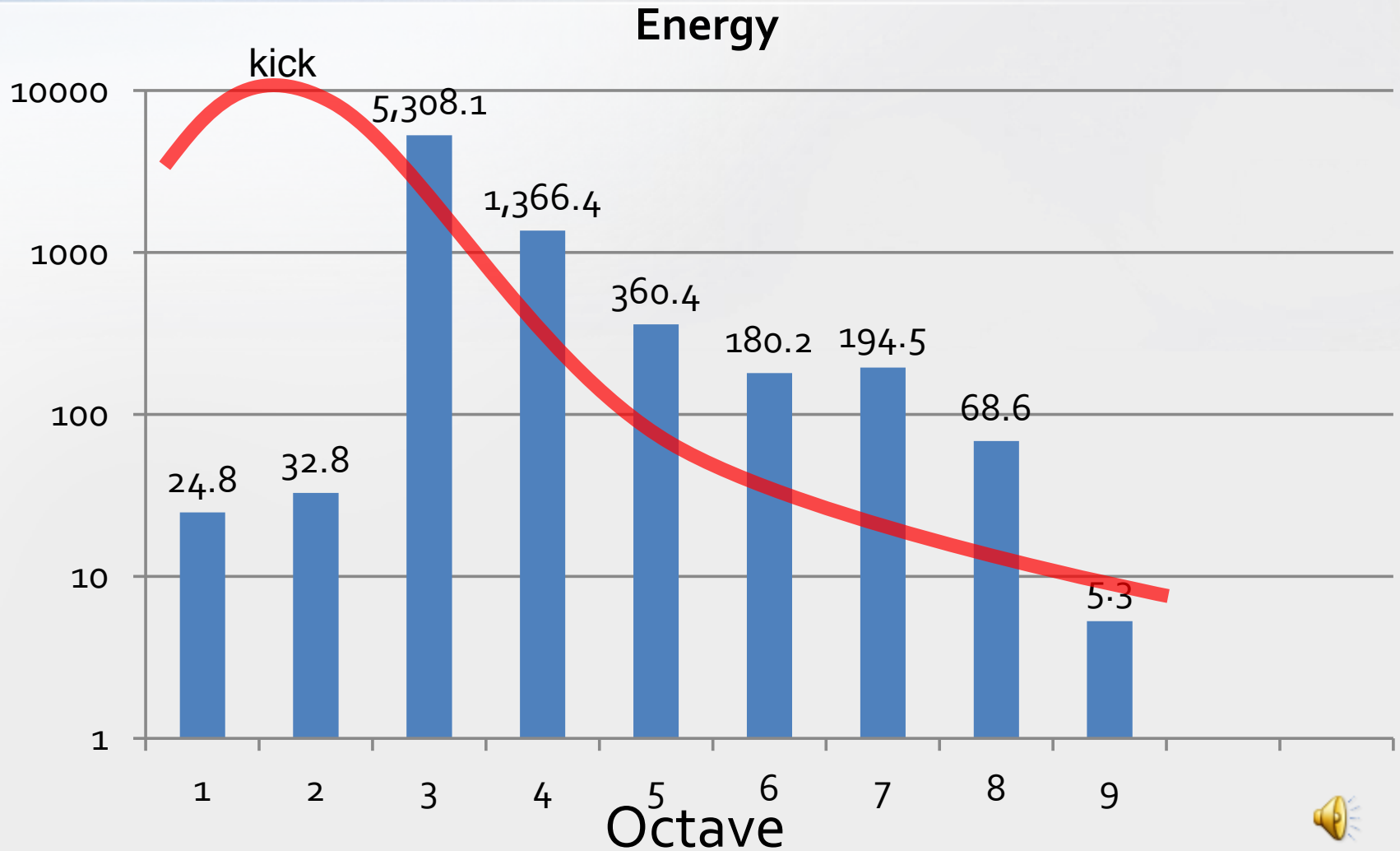# Temporal Information

- Temporal Centroid

# Frame 1

# Features – Frame 1

| Frame | ZCR | Centroid | BW | Skew | Kurtosis | E1 | E2 | E3 | E4 | E5 | E6 | E7 | E8 | E9 |
|-------|-----|----------|-----|------|----------|------|-------|-----|-----|----|----|----|----|----|
| 1 | 9 | 2.8kHz | 5kHz | 2.2 | 6.7 | 4000 | 10100 | 545 | 187 | 77 | 35 | 18 | 9 | 6 |

# Frame 2

# Features : SimpleLoop.wav

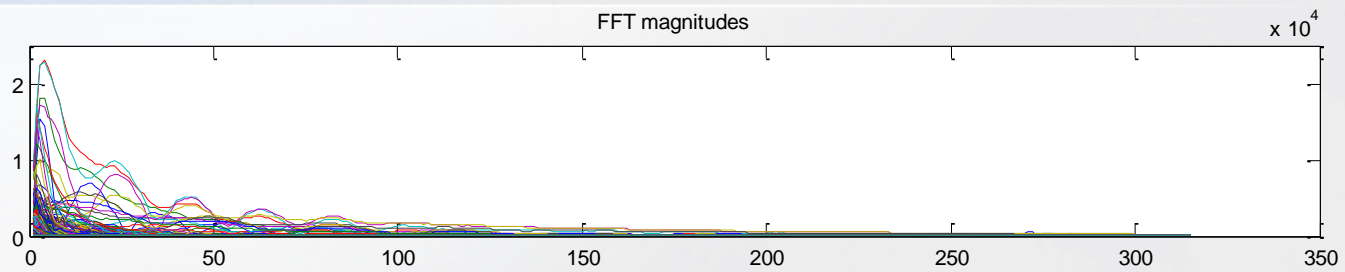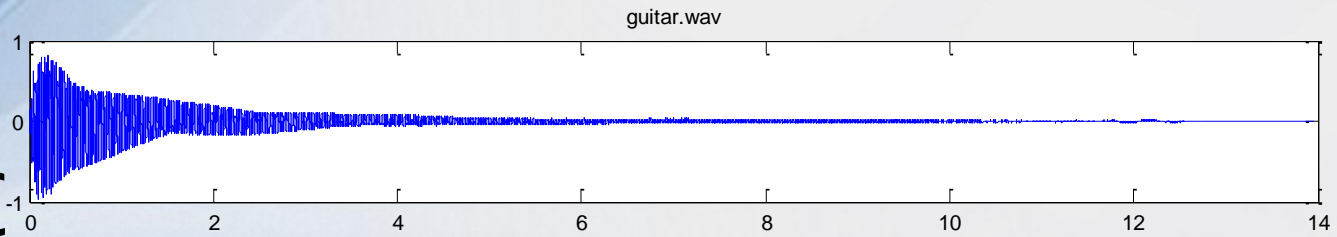| Frame | ZCR | Centroid | BW | Skew | Kurtosis | E1 | E2 | E3 | E4 | E5 | E6 | E7 | E8 |
|-------|-----|----------|------|------|----------|------|-------|------|------|-----|-----|-----|----|
| 1 | 9 | 2.8kHz | 5kHz | 2.2 | 6.7 | 4000 | 10100 | 545 | 187 | 77 | 35 | 18 | 9 |
| 2 | 423 | 3.1kHz | 4kHz | 2 | 7.2 | 24 | 33 | 5300 | 1366 | 360 | 180 | 194 | 68 |

# MFCCs

The idea of MFCCs is to capture spectrum in accordance with human perception.

1. STFT
2. log(STFT)
3. Perform mel-scaling to group and smooth coefficients. (perceptual weighting)
4. Decorrelate with DCT

*[…continued…]*



Weighting for FFT bins to Mel scale



FFT analysis



mel filterbank output

# MFCC
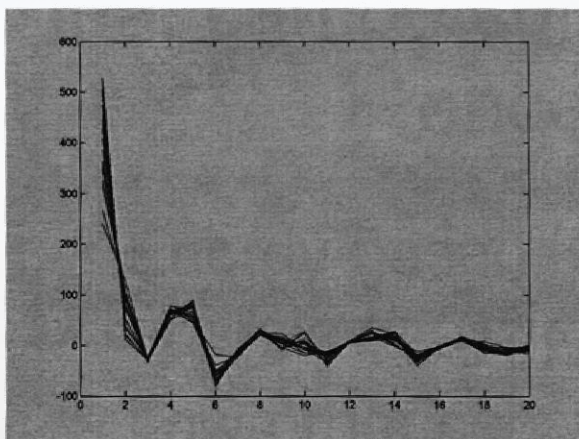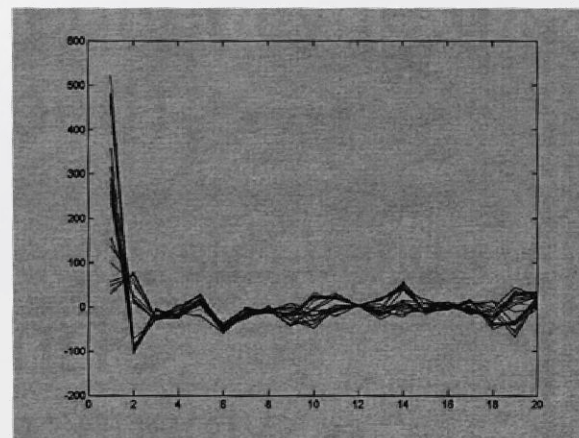


guitar.wav

FFT magnitudes

1

2

3

4

# MFCC of Music
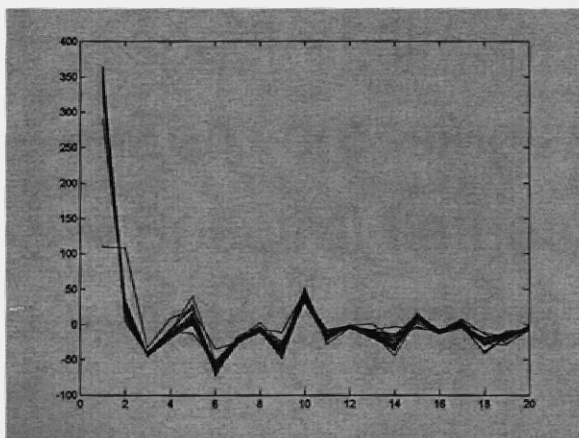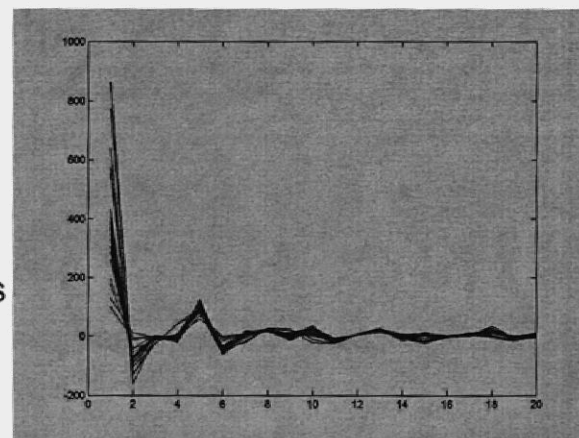
(Petruncio, 2003)

Piano

Saxophone

Tenor
Opera
Singer

Drums

# Features: Measuring changes

- Δ and Δ Δ
  - Change between frames
  - How quickly the change is occurring

- Spectral flux is the distance between the spectrum of successive frames
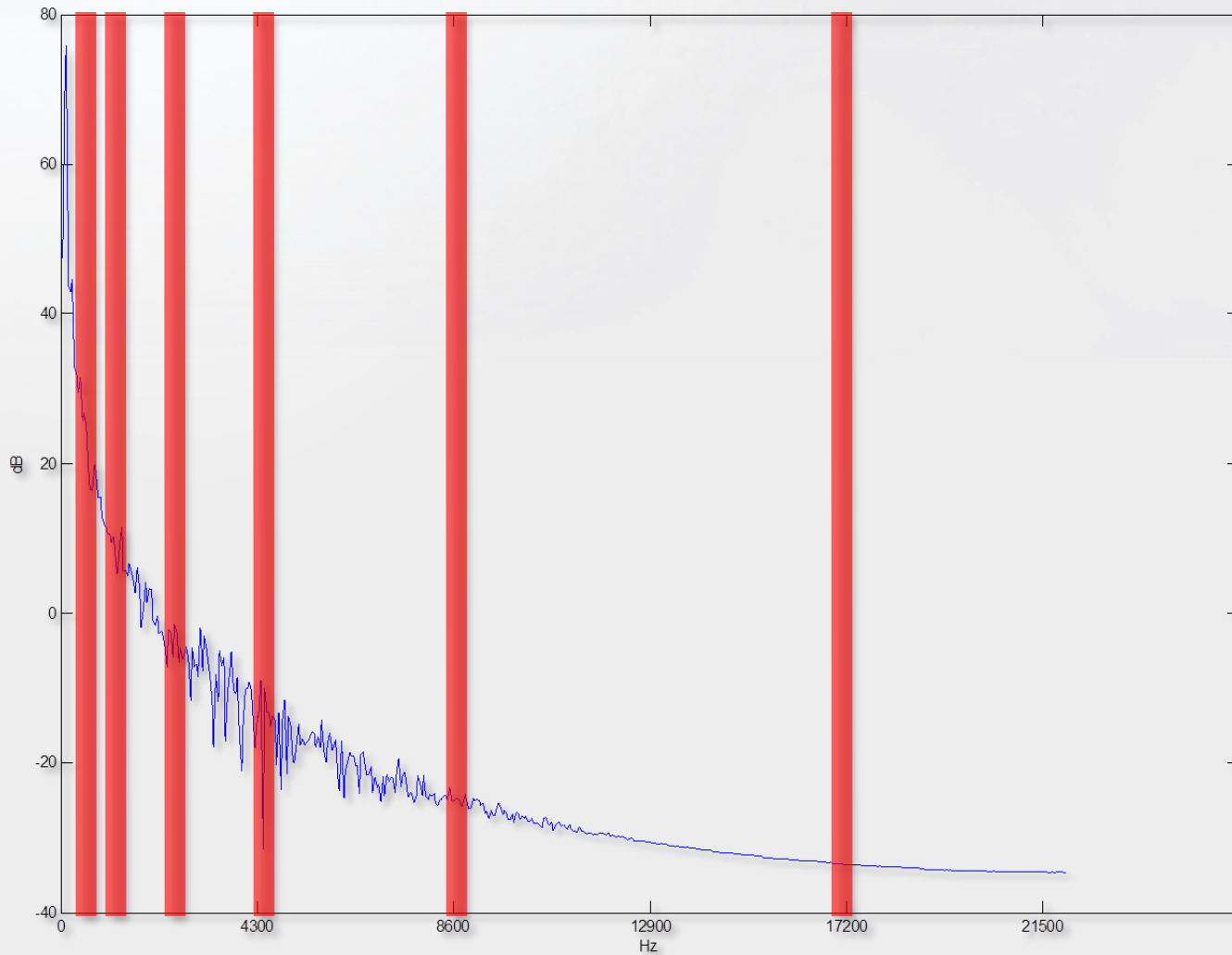
# Spectral Features

- Spectral Flatness Measure
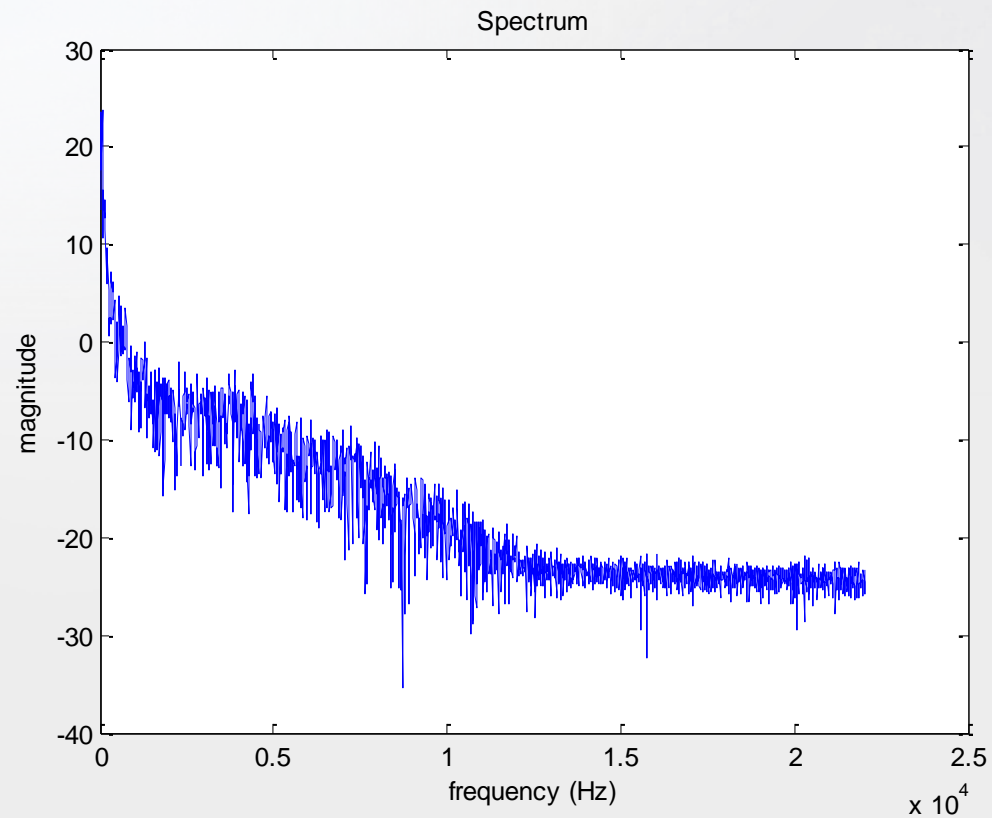- Spectral Crest Factor
- Spectral Flux
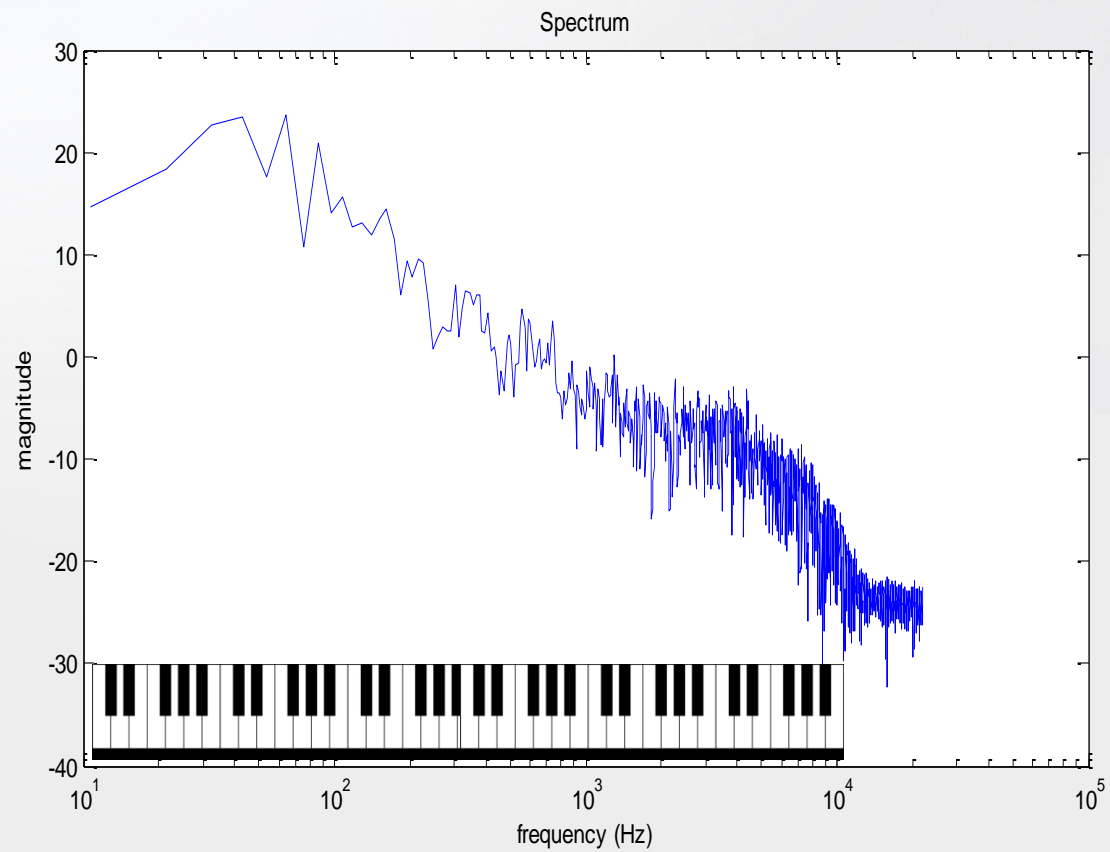
# Feature extraction

- Feature design and creation uses one's domain knowledge.
- Choosing discriminating features is critical
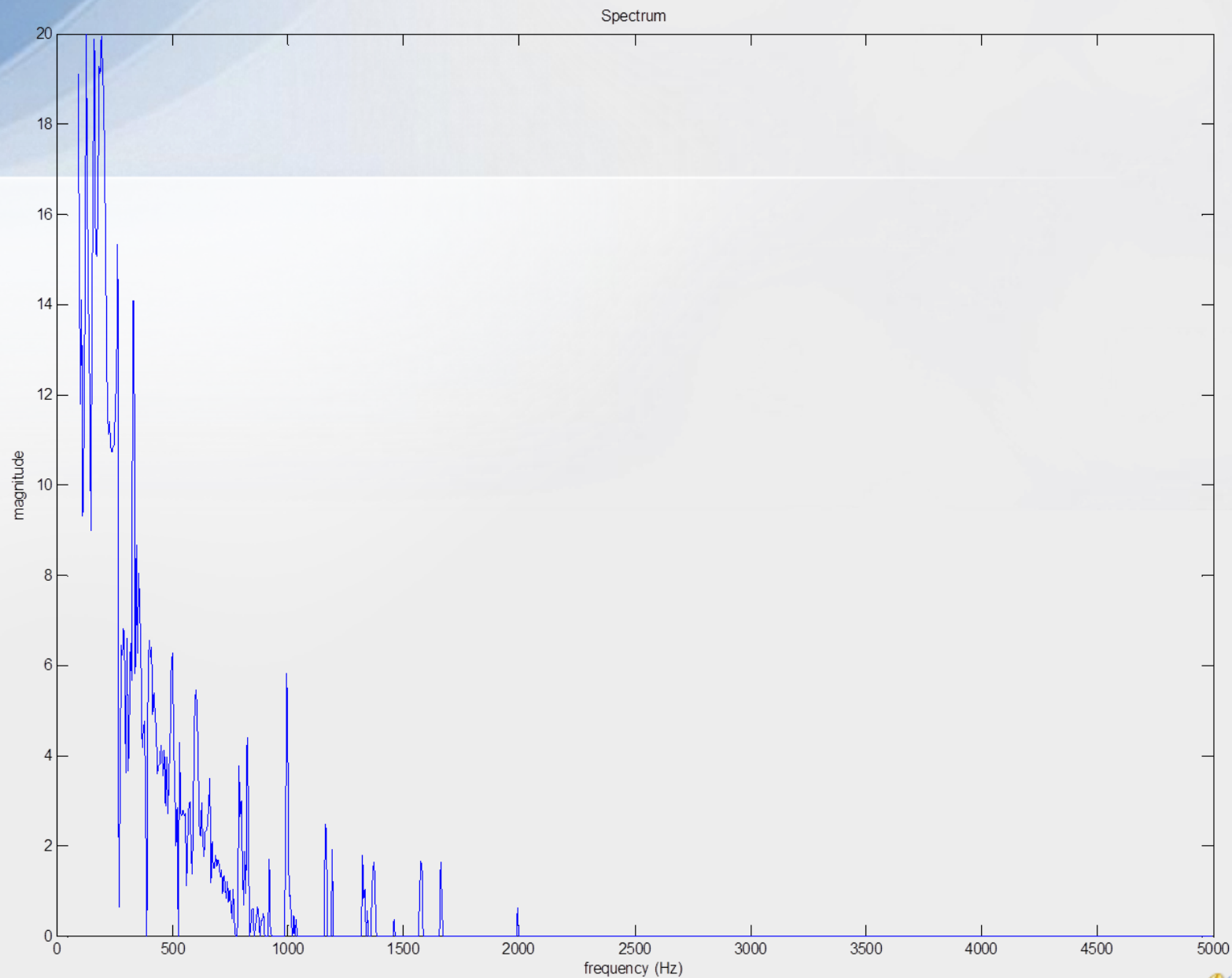- Smaller feature space yields smaller, simpler models, faster training, often less training data needed
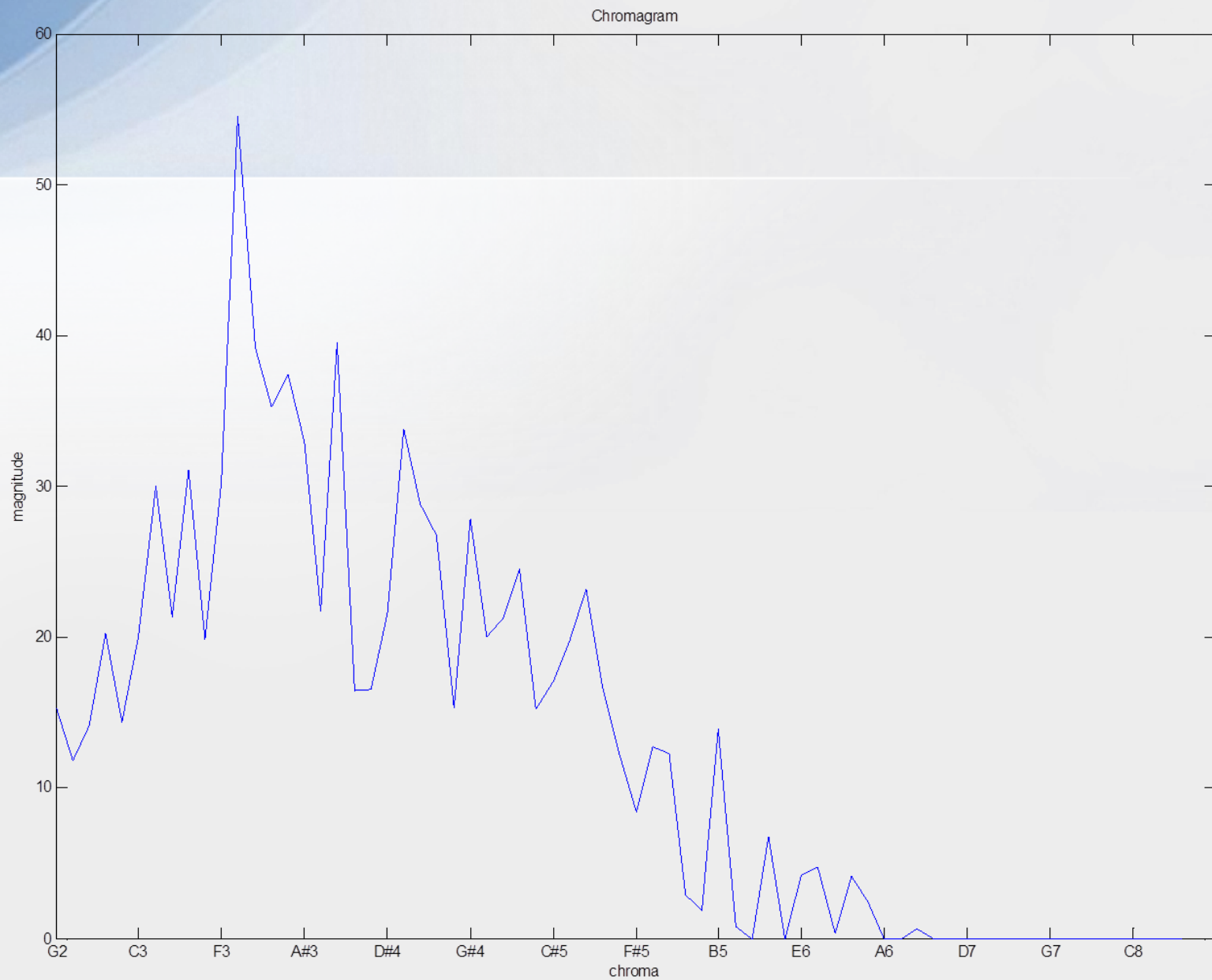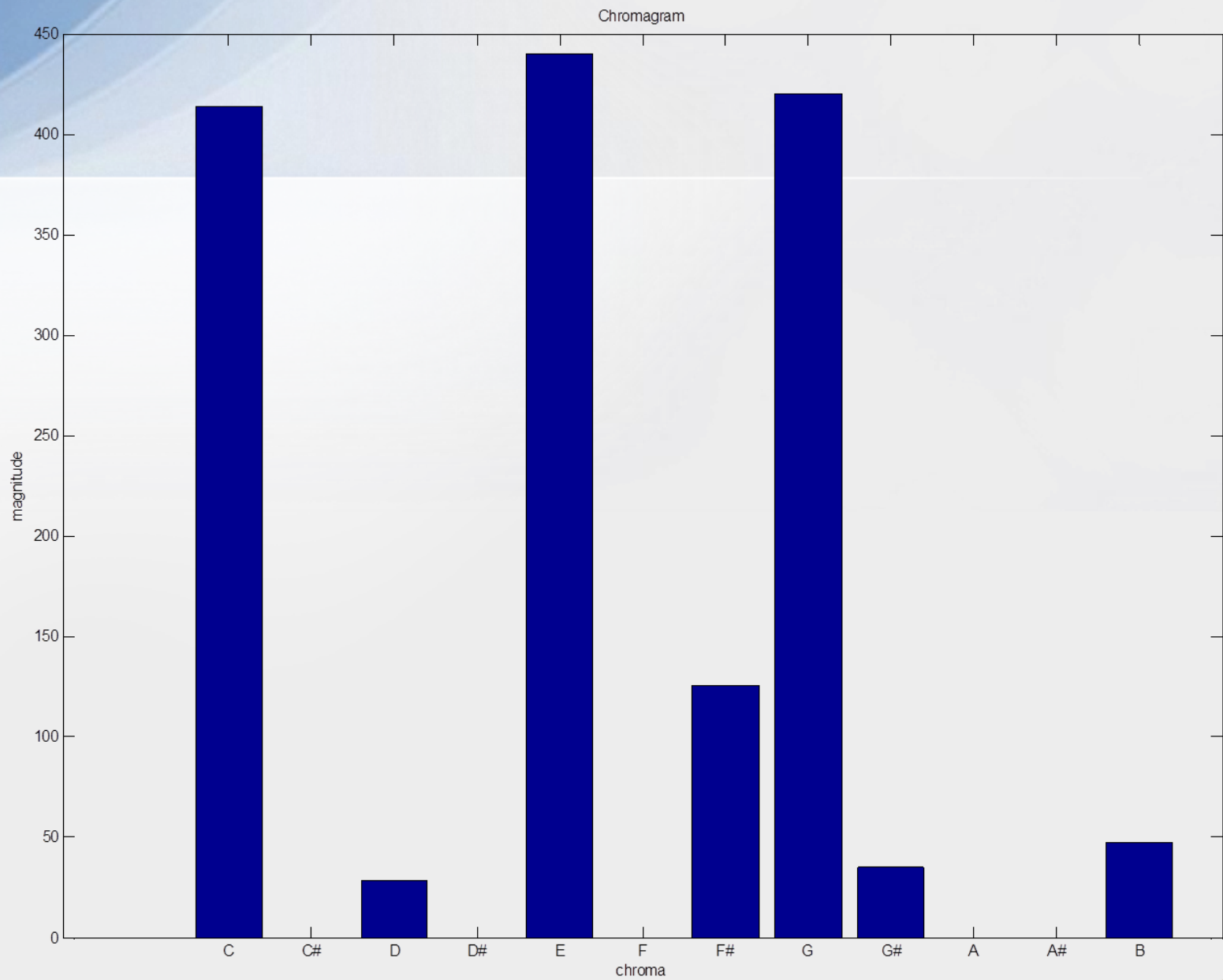
# Spectral Bands
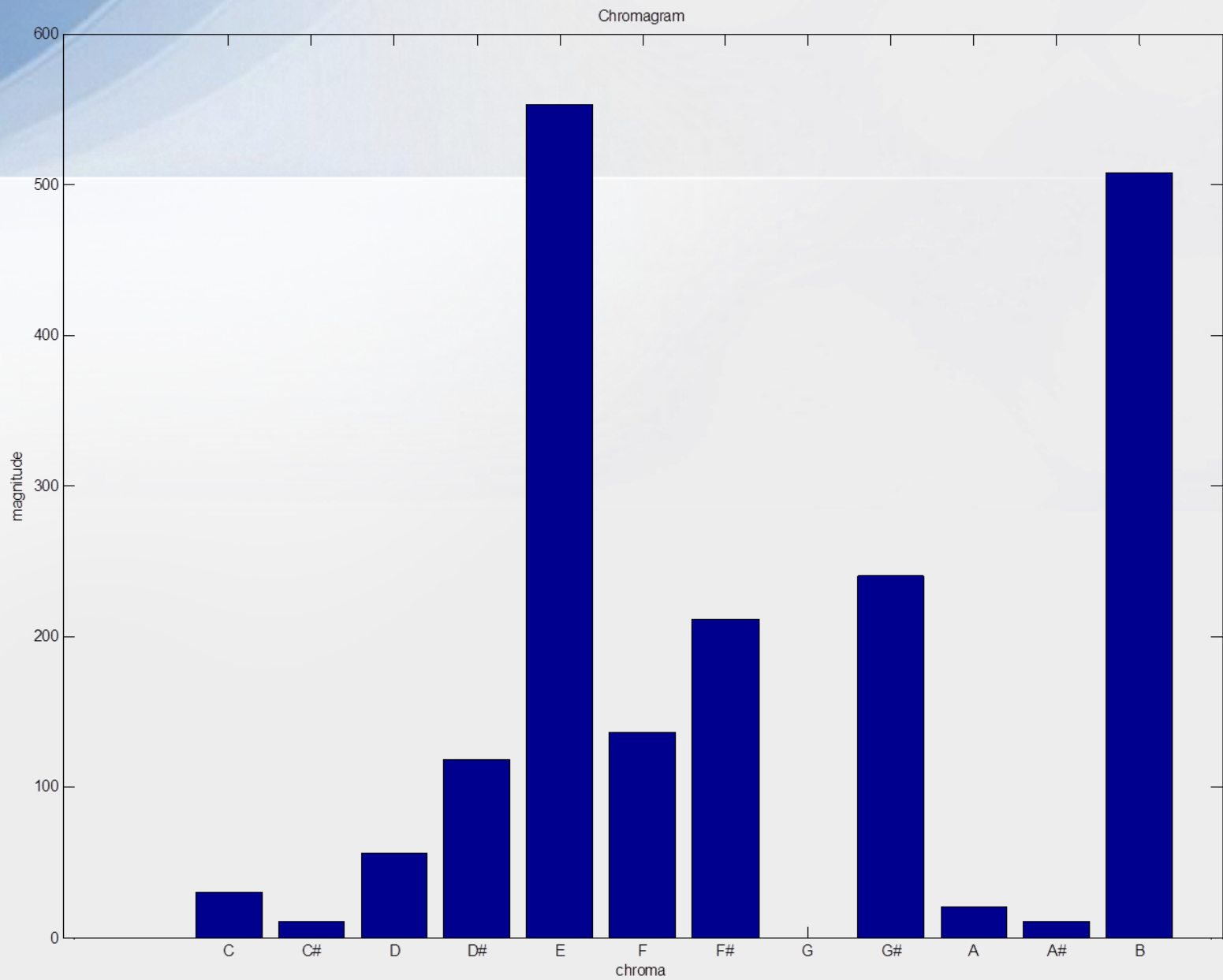
# Log Spectrogram

# Chroma Bins

Spectrum

# EXAMPLE



Picture courtesy: Olivier Lartillot

C major    Cross-Correlations    C minor

C# major    C# minor

D major    D minor

The resulting graph indicate the cross-correlation score for each different tonality candidate.

Key strength

— maj
— min

strength

tonal center

C   C#   D   D#   E   F   F#   G   G#   A   A#   B

- http://www.chordpickout.com/index.html

# Decision stumps

- An example dataset:



This section contains slides adapted from Rob Schapire @ Princeton.

# A decision threshold

- Single threshold: e.g., "output '+' iff x < .2"



- Decision stump: 1 threshold decision

# Many thresholds: Decision trees

- Consists of many decisions in succession (like a flowchart)
- General approach:
  - Recursively split training data into subsets based on simple thresholds
  - Optionally prune to avoid overfitting
- Common algorithms: CART, ID3 => C4.5 (J48)

# Decision Trees

- Advantages:
  - Easy to interpret
  - Decision boundary is explicit and straightforward
- Disadvantages:
  - Can take a long time to learn
    - Finding optimal tree can be NP-complete
  - Prone to overfitting
  - Inherently heuristic
  - Slight perturbations of data can lead to very different trees
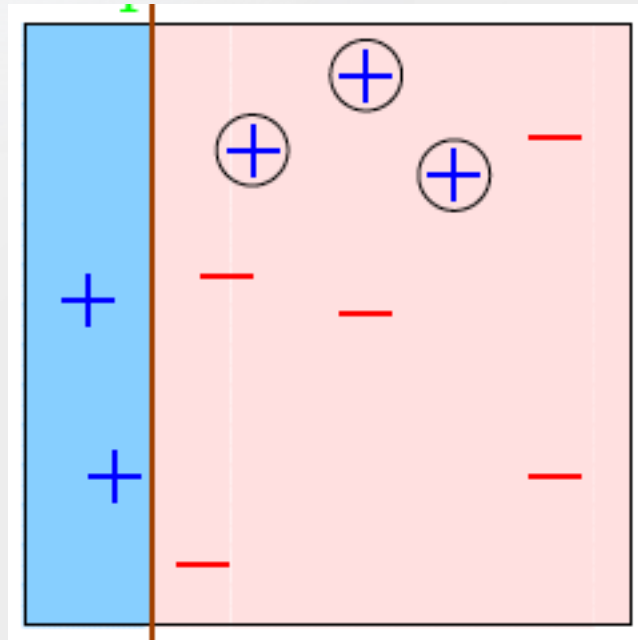
# Boosting

- A "meta-algorithm" for creating a "strong" learner from many "weak" learners
- Iteratively train weak learners on variations of the dataset and combine in a principled way to produce classification outputs.

# AdaBoost

- A popular boosting algorithm from Freund and Schapire
- Robust to overfitting: emphasis on **maximizing the margin**

# Back to stumps
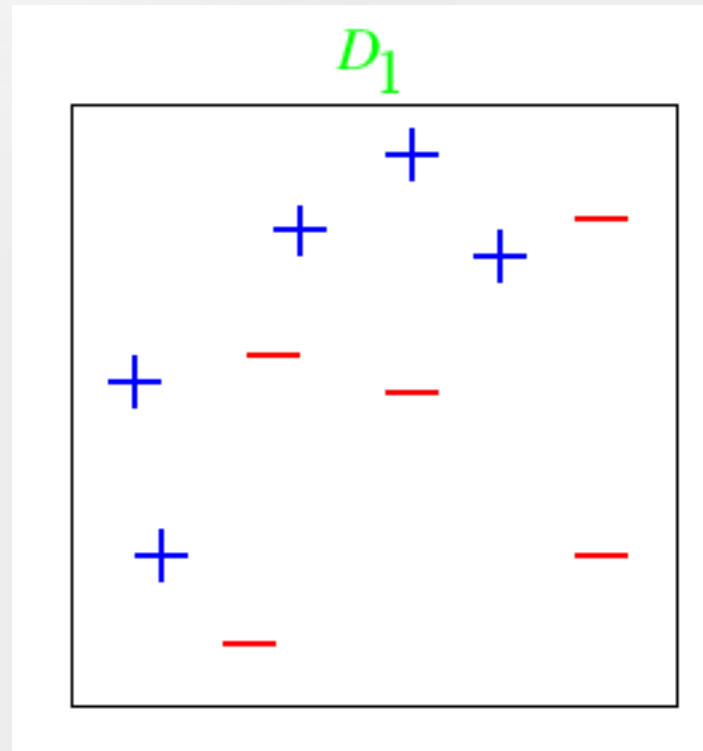
- Single threshold: e.g., "output '+' iff $x < .2$"
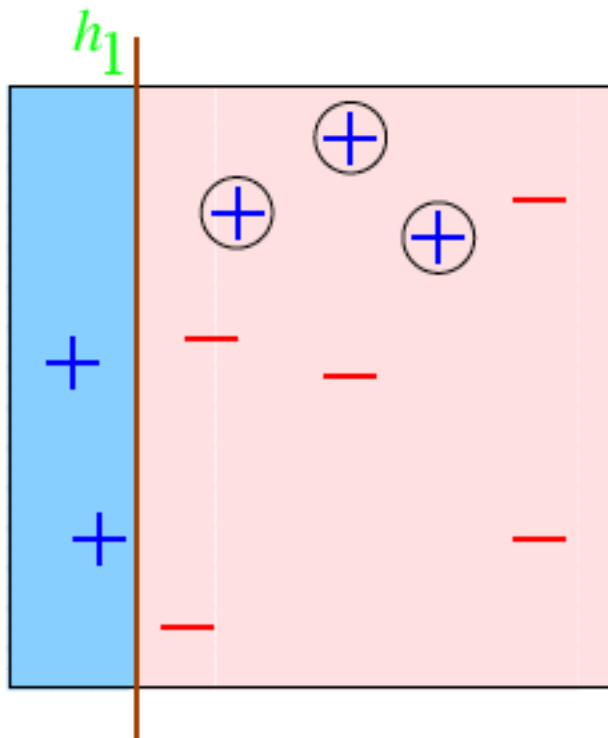


- Makes a nice weak learner!

# The AdaBoost algorithm

- Initialize $D_1$ to be the dataset with each example equally weighted.
- for round t in 1 to T:
  - Train a weak learner, $h_t$, on the dataset $D_t$
  - If $h_t$ can't achieve 50% accuracy, stop.
  - Choose alpha$_t$ according to error rate of ht on $D_t$ (better ht => higher alpha$_t$)
  - Update data weights $D_{t+1}$ to **increase** weight of examples ht got wrong, and **decrease** weight of examples $h_t$ got right.
- To classify new data, take a weighted majority vote of all weak learners, each $h_t$ weighted by its alpha$_t$.
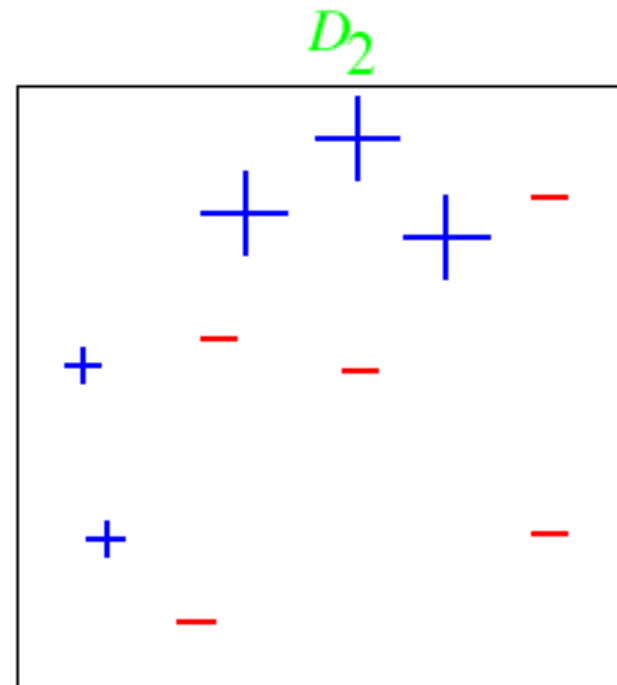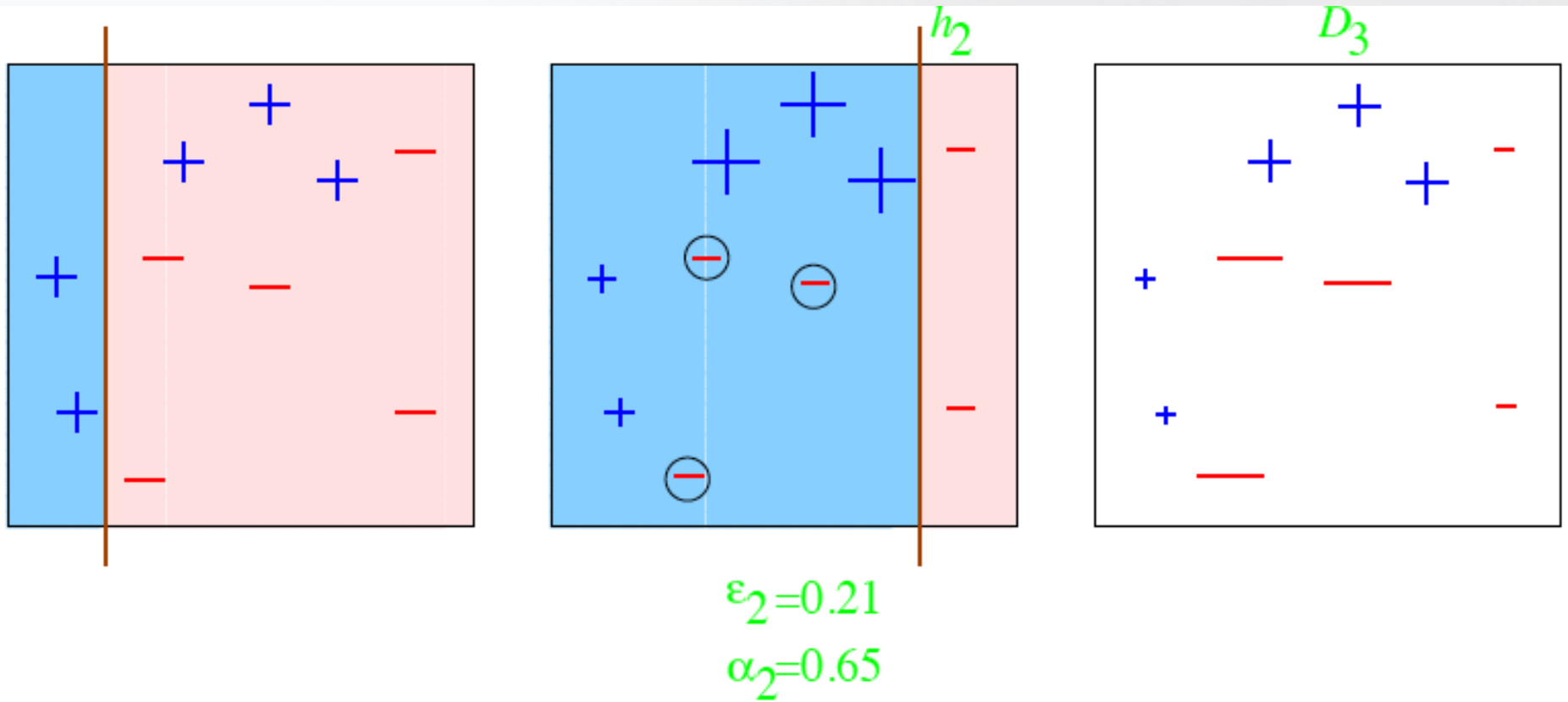
# AdaBoost illustrated

- Initial data:
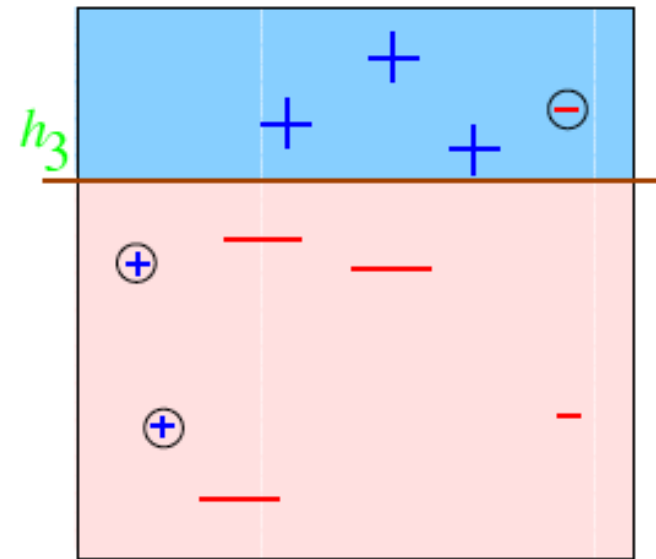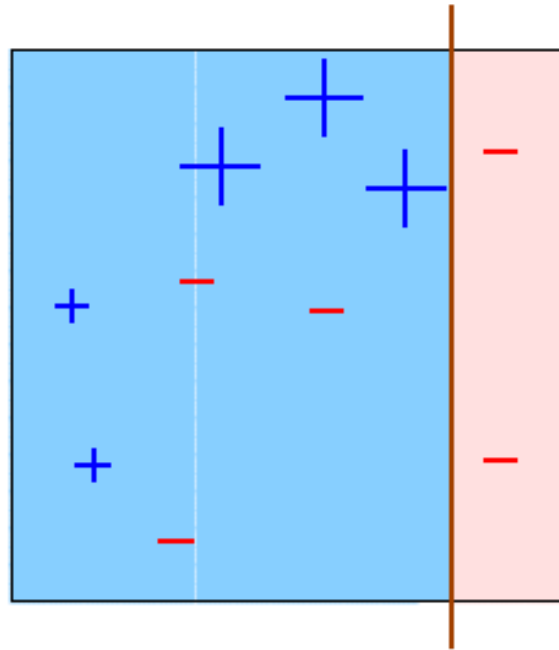
# Round 1



$\varepsilon_1 = 0.30$

$\alpha_1 = 0.42$
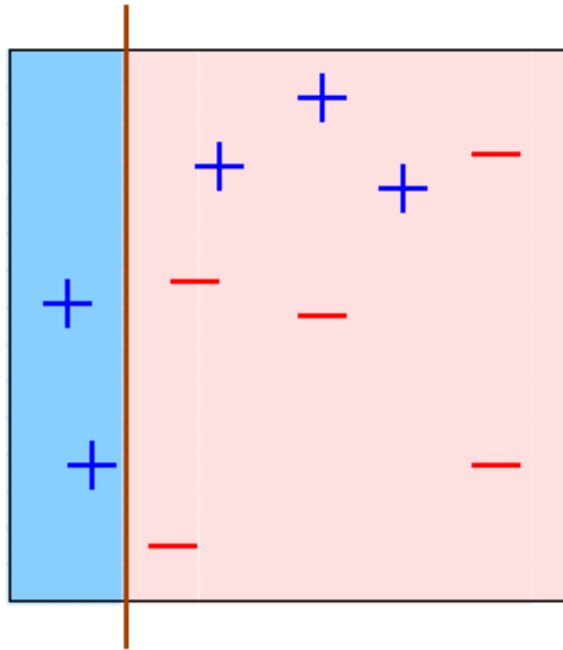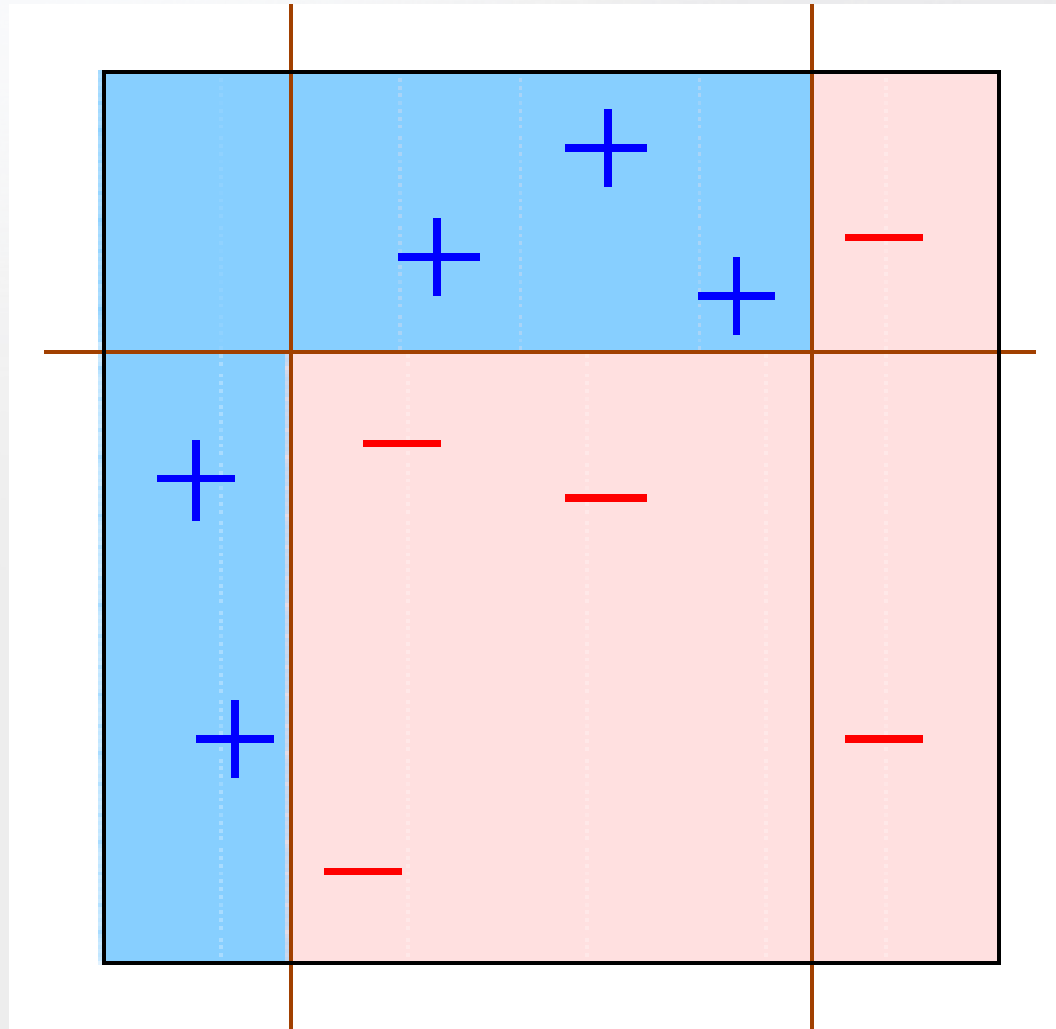
# Round 2

# Round 3



$h_3$

$\varepsilon_3 = 0.14$

$\alpha_3 = 0.92$

# Final classifier
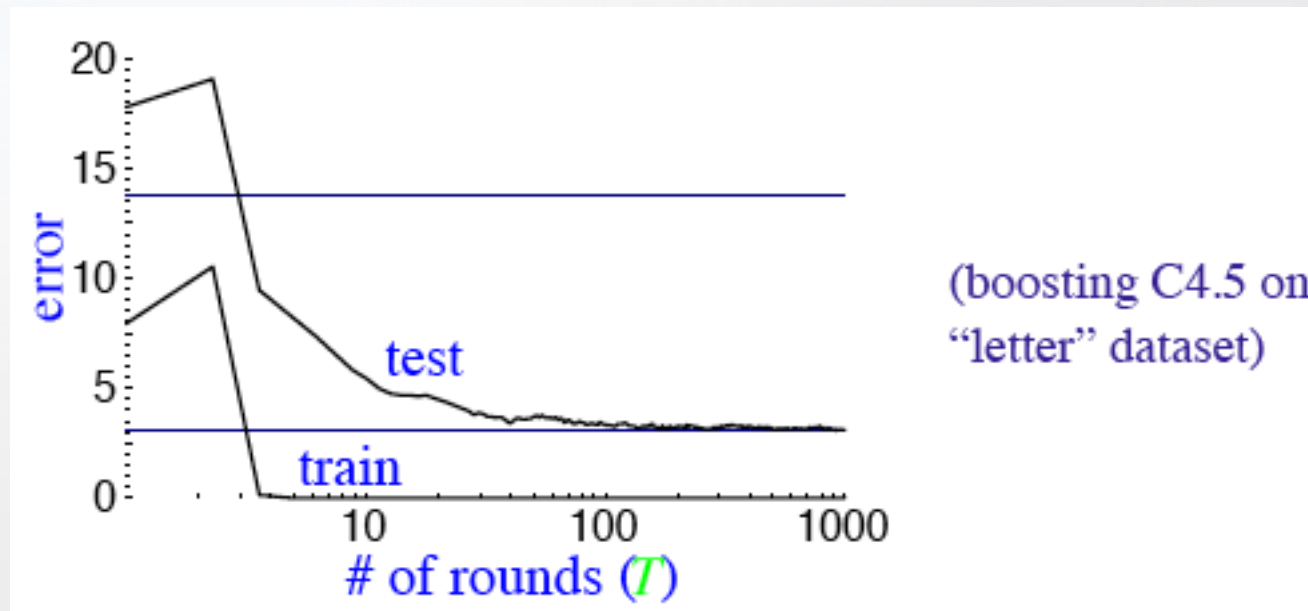


$$H_{\text{final}} = \text{sign}\left( 0.42 \quad + 0.65 \quad + 0.92 \right)$$

# Final classifier: decision boundary

# A typical AdaBoost run



(boosting C4.5 on "letter" dataset)

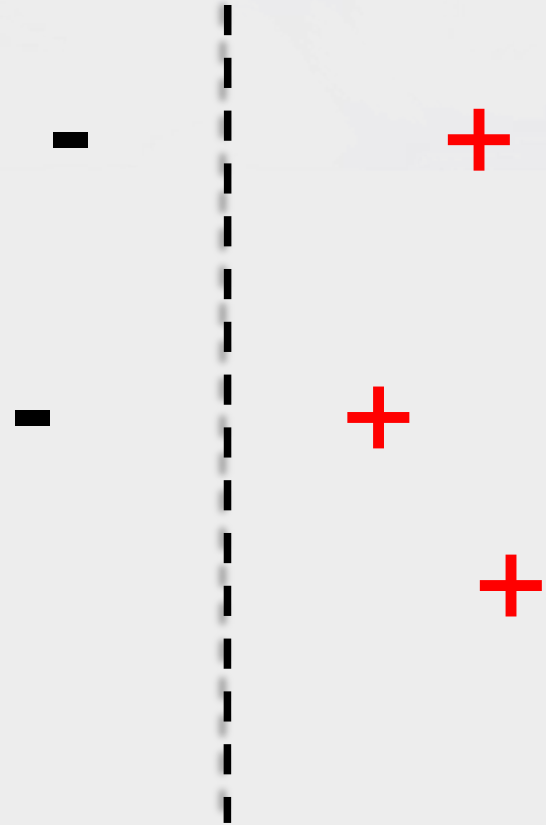- Test error does not increase, even after 1000 rounds
- Test error continues to drop, even after training error = 0.

# The margin

- Narrow margin
- Wide margin

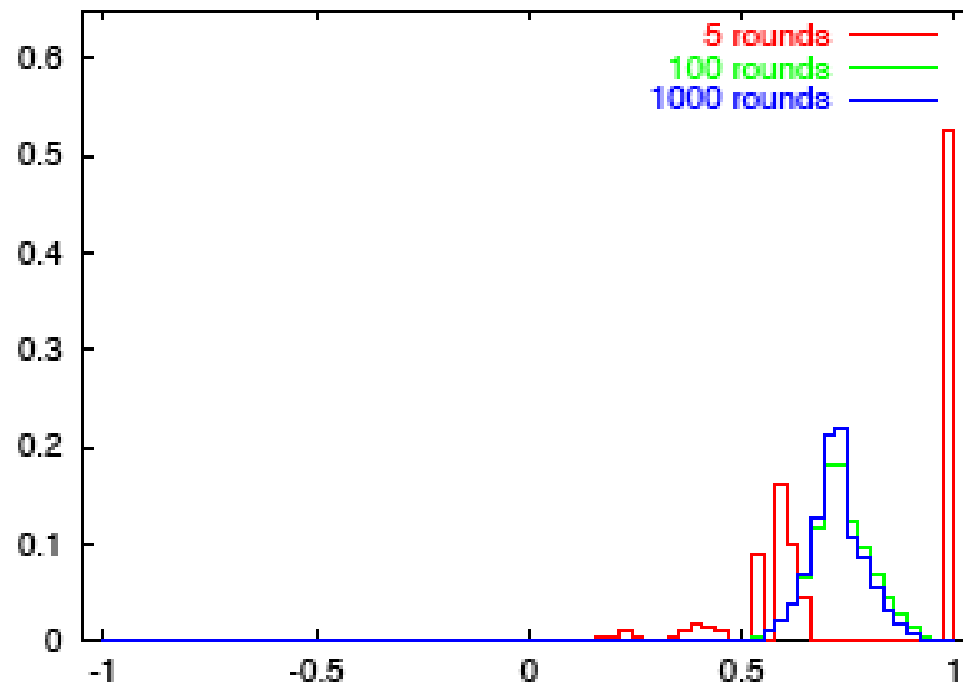# Margin distribution after N rounds



| | # rounds | | |
|---|---|---|---|
| | 5 | 100 | 1000 |
| train error | 0.0 | 0.0 | 0.0 |
| test error | 8.4 | 3.3 | 3.1 |
| % margins $\leq 0.5$ | 7.7 | 0.0 | 0.0 |
| minimum margin | 0.14 | 0.52 | 0.55 |

# AdaBoost pro & con

- Advantages:
  - Robust to overfitting
  - Conceptually simple
  - Statistically very nice: maximizing the margin, game-theoretic understanding
  - Can work with any base learner
  - No parameters to tune
- Disadvantages:
  - Weak learner must achieve >50% or failure
  - Original formulation binary only
    - AdaBoost.M1 handles multi-class, but more required of weak learner
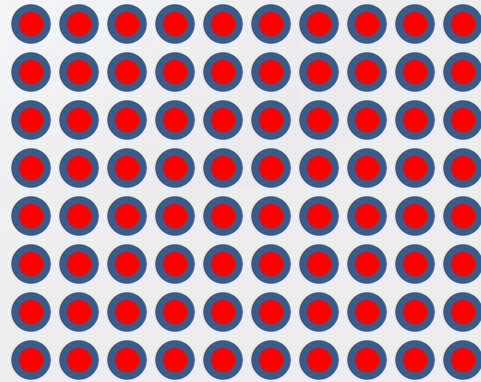
# EVALUATION
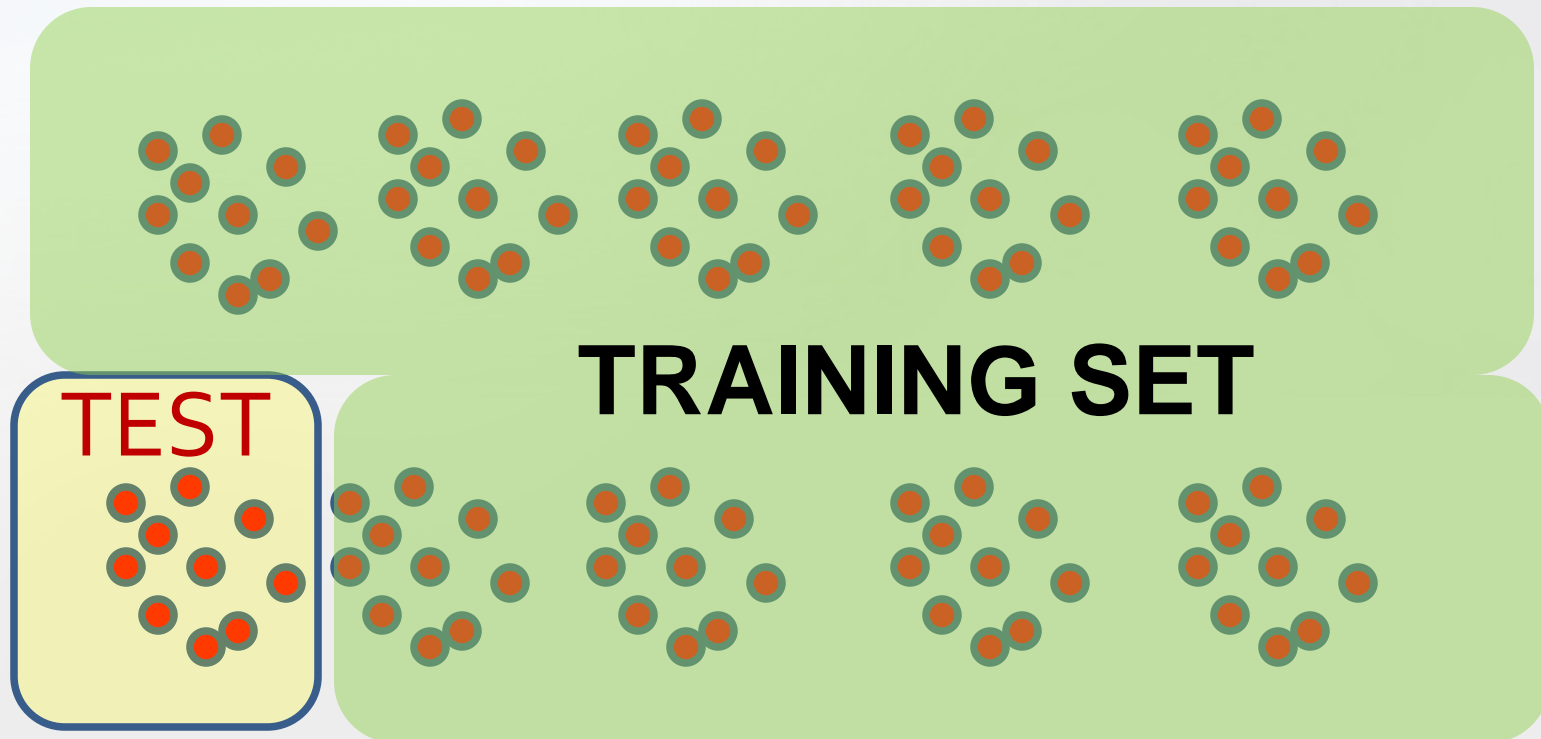
Our classifier accuracy is 83.4%

# Cross-validation

- Say, 10-fold cross validation
- Divide test set into 10 random subsets.
- 1 test set is tested using the classifier trained on the remaining 9.
- We then do test/train on all of the other sets and average the percentages. Helps prevent over fitting.
- Do not optimize too much on cross validation – you can severely overfit. Sanity check with a test set.
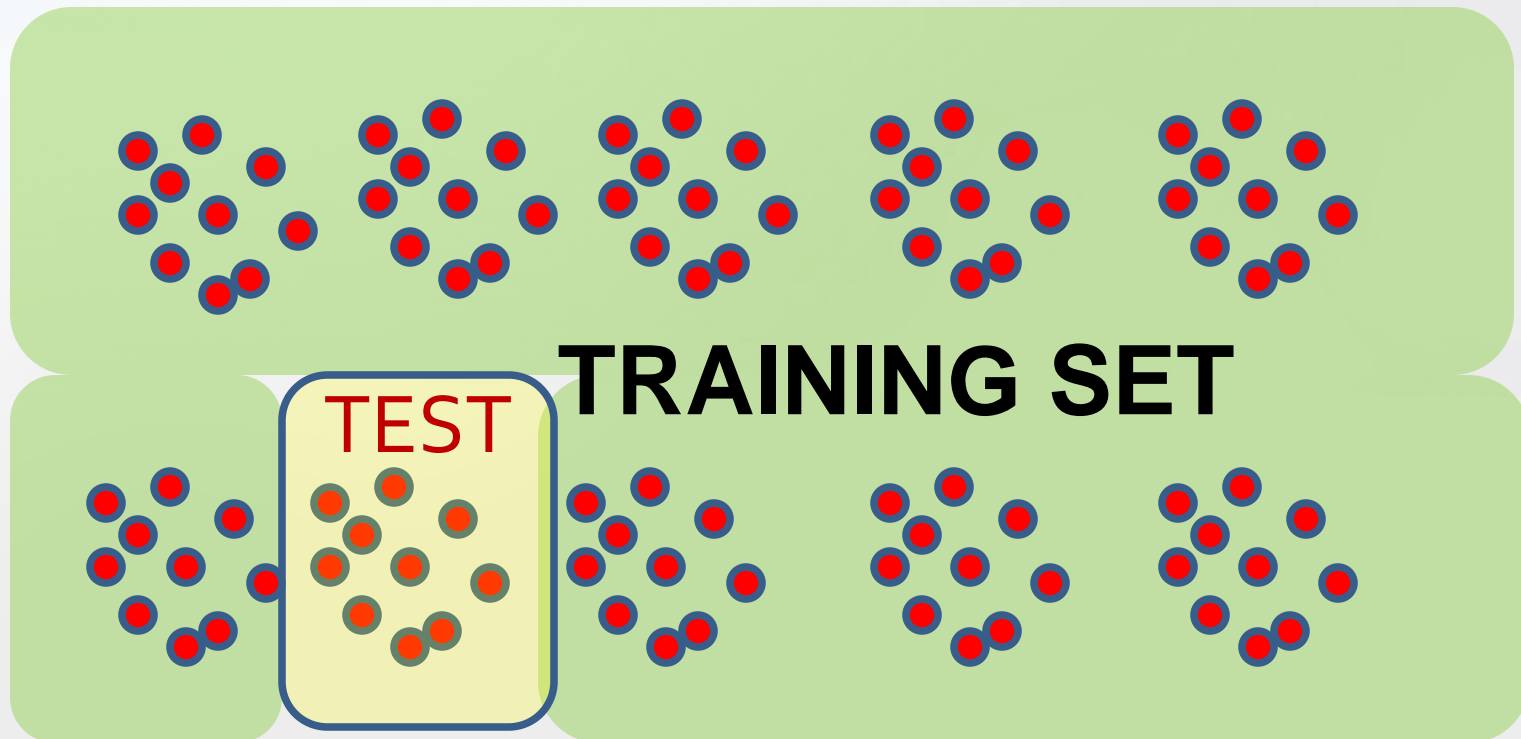
# Cross-validation

# Cross-validation



**TRAINING SET**

TEST

Fold 1: 70%

# Cross-validation



TRAINING SET

TEST

Fold 1: 70%
Fold 2: 80%

# Cross-validation

Fold 1: 76%
Fold 2: 80%
Fold 3: 77%
Fold 4: 83%
Fold 5: 72%
Fold 6: 82%
Fold 7: 81%
Fold 8: 71%
Fold 9: 90%
Fold 10: 82%
**Mean = 79.4%**

# Stratified Cross-Validation

- Same as cross-validation, except that the folds are chosen so that they contain equal proportions of labels.

# > End Day 2