

Analyzing training error

Analysis reveals:

- What α_t to choose for hypothesis h_t ?
What h_t to choose?

$$\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$$

ϵ_t - weighted training error

- If each weak learner h_t is slightly better than random guessing ($\epsilon_t < 0.5$), then training error of AdaBoost decays exponentially fast in number of rounds T .

$$\frac{1}{m} \sum_{i=1}^m \delta(H(x_i) \neq y_i) \leq \exp \left(-2 \sum_{t=1}^T (1/2 - \epsilon_t)^2 \right)$$

Training Error

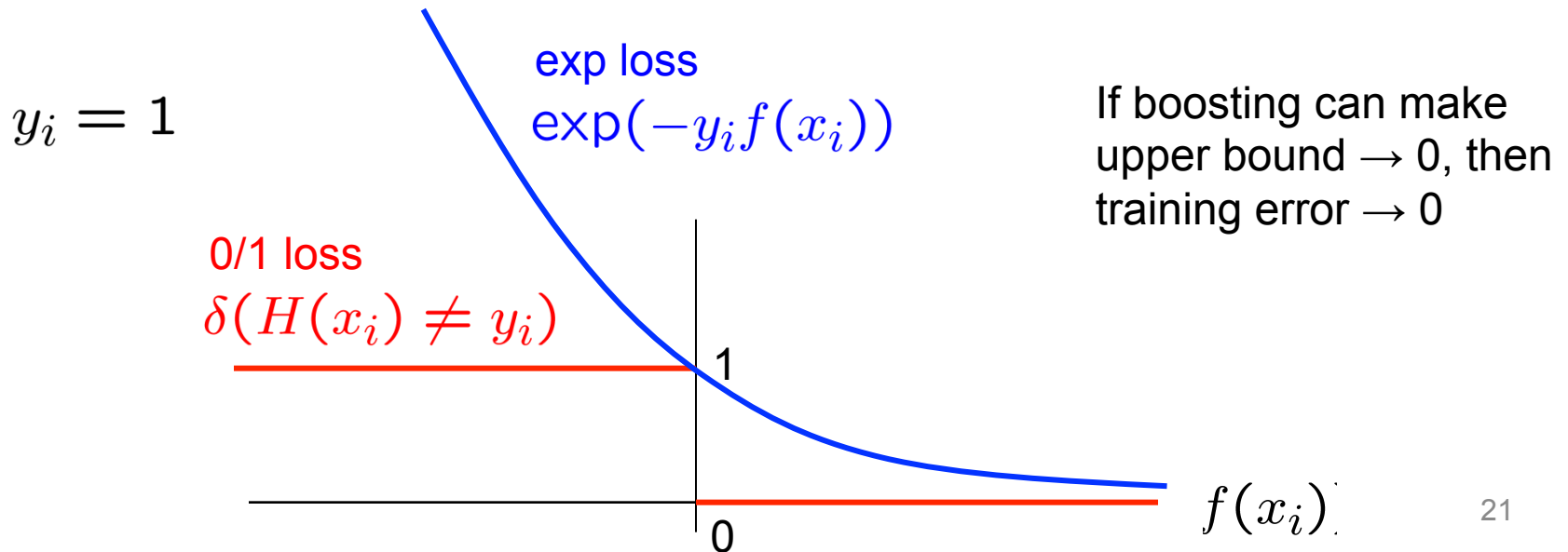
Analyzing training error

Training error of final classifier is bounded by:

$$\frac{1}{m} \sum_{i=1}^m \delta(H(x_i) \neq y_i) \leq \frac{1}{m} \sum_{i=1}^m \exp(-y_i f(x_i))$$

Convex
upper
bound

Where $f(x) = \sum_t \alpha_t h_t(x)$; $H(x) = \text{sign}(f(x))$



Analyzing training error

Training error of final classifier is bounded by:

$$\frac{1}{m} \sum_{i=1}^m \delta(H(x_i) \neq y_i) \leq \frac{1}{m} \sum_{i=1}^m \exp(-y_i f(x_i)) = \prod_t Z_t$$

Where $f(x) = \sum_t \alpha_t h_t(x)$; $H(x) = \text{sign}(f(x))$

Proof:

Using Weight Update Rule

$$D_1(i) = 1/m.$$

$$D_2(i) = \frac{1}{m} \frac{e^{-\alpha_1 y_i h_1(x_i)}}{Z_1}$$

$$D_3(i) = \frac{1}{m} \frac{e^{-\alpha_1 y_i h_1(x_i)} e^{-\alpha_2 y_i h_2(x_i)}}{Z_1 Z_2}$$

...

$$D_{T+1}(i) = \frac{1}{m} \frac{\exp(-y_i f(x_i))}{\prod_t Z_t}$$

Wts of all pts add to 1

$$\sum_{i=1}^m D_{T+1}(i) = 1$$

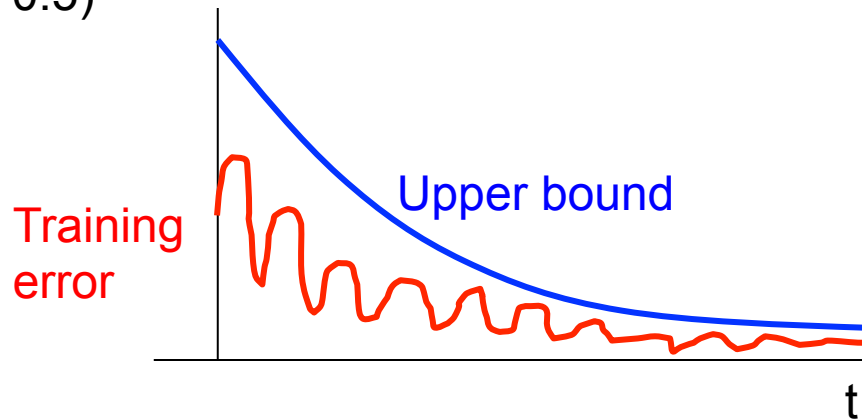
Analyzing training error

Training error of final classifier is bounded by:

$$\frac{1}{m} \sum_{i=1}^m \delta(H(x_i) \neq y_i) \leq \frac{1}{m} \sum_{i=1}^m \exp(-y_i f(x_i)) = \prod_t Z_t$$

Where $f(x) = \sum_t \alpha_t h_t(x)$; $H(x) = \text{sign}(f(x))$

If $Z_t < 1$, training error decreases exponentially (even though weak learners may not be good $\varepsilon_t \sim 0.5$)



What α_t to choose for hypothesis h_t ?

Training error of final classifier is bounded by:

$$\frac{1}{m} \sum_{i=1}^m \delta(H(x_i) \neq y_i) \leq \frac{1}{m} \sum_{i=1}^m \exp(-y_i f(x_i)) = \prod_t Z_t$$

Where $f(x) = \sum_t \alpha_t h_t(x)$; $H(x) = \text{sign}(f(x))$

If we minimize $\prod_t Z_t$, we minimize our training error

We can tighten this bound greedily, by choosing α_t and h_t on each iteration to minimize Z_t .

$$Z_t = \sum_{i=1}^m D_t(i) \exp(-\alpha_t y_i h_t(x_i))$$

What α_t to choose for hypothesis h_t ?

We can minimize this bound by choosing α_t on each iteration to minimize Z_t .

$$Z_t = \sum_{i=1}^m D_t(i) \exp(-\alpha_t y_i h_t(x_i))$$

For boolean target function, this is accomplished by [Freund & Schapire '97]:

$$\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$$

Proof:

$$\begin{aligned} Z_t &= \sum_{i: y_i \neq h_t(x_i)} D_t(i) e^{\alpha_t} + \sum_{i: y_i = h_t(x_i)} D_t(i) e^{-\alpha_t} \\ &= \epsilon_t e^{\alpha_t} + (1 - \epsilon_t) e^{-\alpha_t} \end{aligned}$$

$$\frac{\partial Z_t}{\partial \alpha_t} = \epsilon_t e^{\alpha_t} - (1 - \epsilon_t) e^{-\alpha_t} = 0 \quad \Rightarrow \quad e^{2\alpha_t} = \frac{1 - \epsilon_t}{\epsilon_t}$$

Bounding the error with choice of α_t

We can minimize this bound by choosing α_t on each iteration to minimize Z_t .

$$Z_t = \sum_{i=1}^m D_t(i) \exp(-\alpha_t y_i h_t(x_i))$$

For boolean target function, this is accomplished by [Freund & Schapire '97]:

$$\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$$

Proof:

$$\begin{aligned} Z_t &= \sum_{i: y_i \neq h_t(x_i)} D_t(i) e^{\alpha_t} + \sum_{i: y_i = h_t(x_i)} D_t(i) e^{-\alpha_t} \\ &= \epsilon_t e^{\alpha_t} + (1 - \epsilon_t) e^{-\alpha_t} \\ &= 2\sqrt{\epsilon_t(1 - \epsilon_t)} = \sqrt{1 - (1 - 2\epsilon_t)^2} \end{aligned}$$

Dumb classifiers made Smart

Training error of final classifier is bounded by:

$$\frac{1}{m} \sum_{i=1}^m \delta(H(x_i) \neq y_i) \leq \prod_t Z_t = \prod_t \sqrt{1 - (1 - 2\epsilon_t)^2}$$

$$\leq \exp \left(-2 \sum_{t=1}^T \underbrace{(1/2 - \epsilon_t)^2}_{\text{grows as } \epsilon_t \text{ moves away from } 1/2} \right)$$

grows as ϵ_t moves
away from 1/2

If each classifier is (at least slightly) better than random $\epsilon_t < 0.5$

AdaBoost will achieve zero training error exponentially fast (in number of rounds T) !!

What about test error?

Boosting and Logistic Regression

Logistic regression assumes:

$$P(Y = 1|X) = \frac{1}{1 + \exp(f(x))}$$

$$f(x) = w_0 + \sum_j w_j x_j$$

And tries to maximize data likelihood:

$$P(\mathcal{D}|f) \stackrel{\text{iid}}{=} \prod_{i=1}^m \frac{1}{1 + \exp(-y_i f(x_i))}$$

Equivalent to minimizing log loss

$$-\log P(\mathcal{D}|f) = \sum_{i=1}^m \ln(1 + \exp(-y_i f(x_i)))$$

Boosting and Logistic Regression

Logistic regression equivalent to minimizing log loss

$$\sum_{i=1}^m \ln(1 + \exp(-y_i f(x_i)))$$

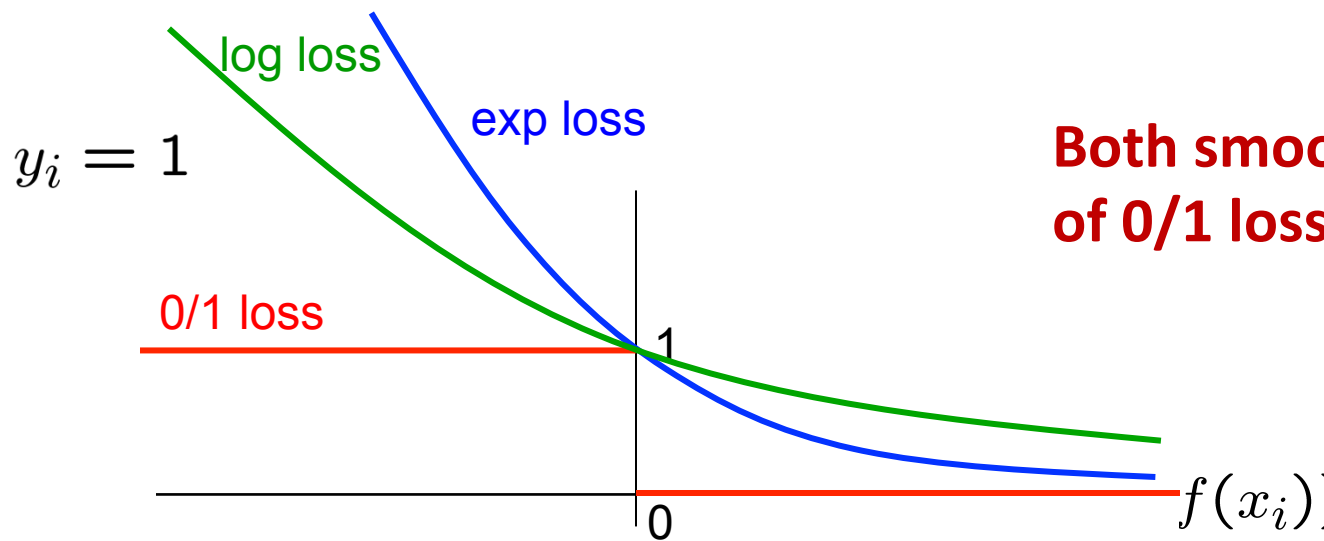
$$f(x) = w_0 + \sum_j w_j x_j$$

Boosting minimizes similar loss function!!

$$\frac{1}{m} \sum_{i=1}^m \exp(-y_i f(x_i))$$

$$f(x) = \sum_t \alpha_t h_t(x)$$

Weighted average of weak learners



Boosting and Logistic Regression

Logistic regression:

- Minimize log loss

$$\sum_{i=1}^m \ln(1 + \exp(-y_i f(x_i)))$$

- Define

$$f(x) = \sum_j w_j x_j$$

where x_j predefined features

(linear classifier)

- Jointly optimize over all weights $w_0, w_1, w_2...$

Boosting:

- Minimize exp loss

$$\sum_{i=1}^m \exp(-y_i f(x_i))$$

- Define

$$f(x) = \sum_t \alpha_t h_t(x)$$

where $h_t(x)$ defined dynamically to fit data

(not a linear classifier)

- Weights α_t learned per iteration t incrementally

Hard & Soft Decision

Weighted average of weak learners $f(x) = \sum_t \alpha_t h_t(x)$

Hard Decision/Predicted label: $H(x) = \text{sign}(f(x))$

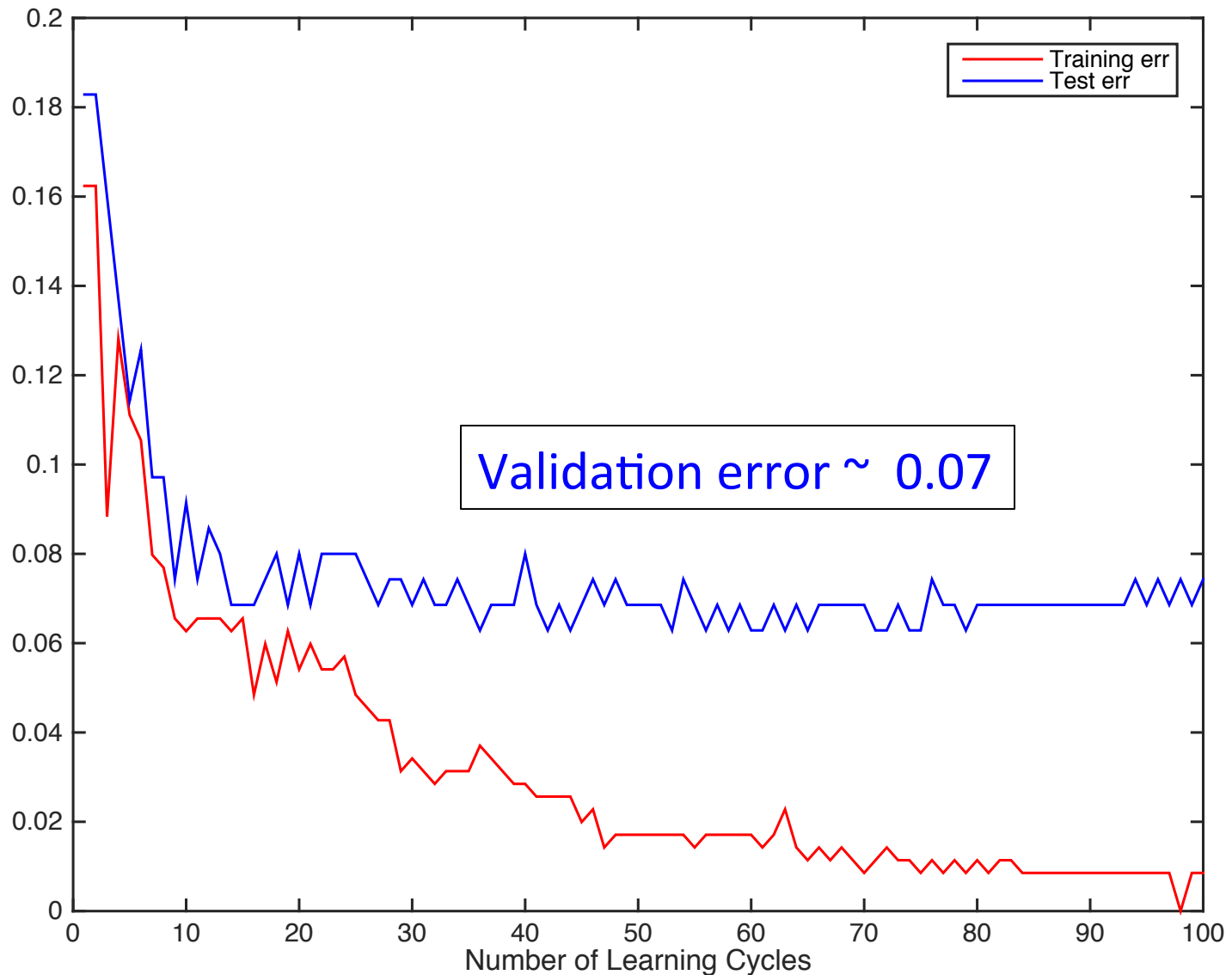
Soft Decision:
(based on analogy with
logistic regression)

$$P(Y = 1|X) = \frac{1}{1 + \exp(f(x))}$$

Matlab example - boosting

- % UCI dataset
- % 34 features, 351 samples
- % binary classification
- load ionosphere;
- rng(2); % For reproducibility
- ClassTreeEns = fitensemble(X,Y,'AdaBoostM1',100,'Tree');
- rsLoss = resubLoss(ClassTreeEns,'Mode','Cumulative');
- plot(rsLoss,'r');
- hold on
- ClassTreeEns = fitensemble(X,Y,'AdaBoostM1',100,'Tree',...
- 'Holdout',0.5);
- genError = kfoldLoss(ClassTreeEns,'Mode','Cumulative');
- plot(genError,'b');
- xlabel('Number of Learning Cycles');
- legend('Training err', 'Test err')

Matlab example - boosting



Matlab example – decision tree

load ionosphere

% UCI dataset

% 34 features, 351 samples

% binary classification

rng(100)

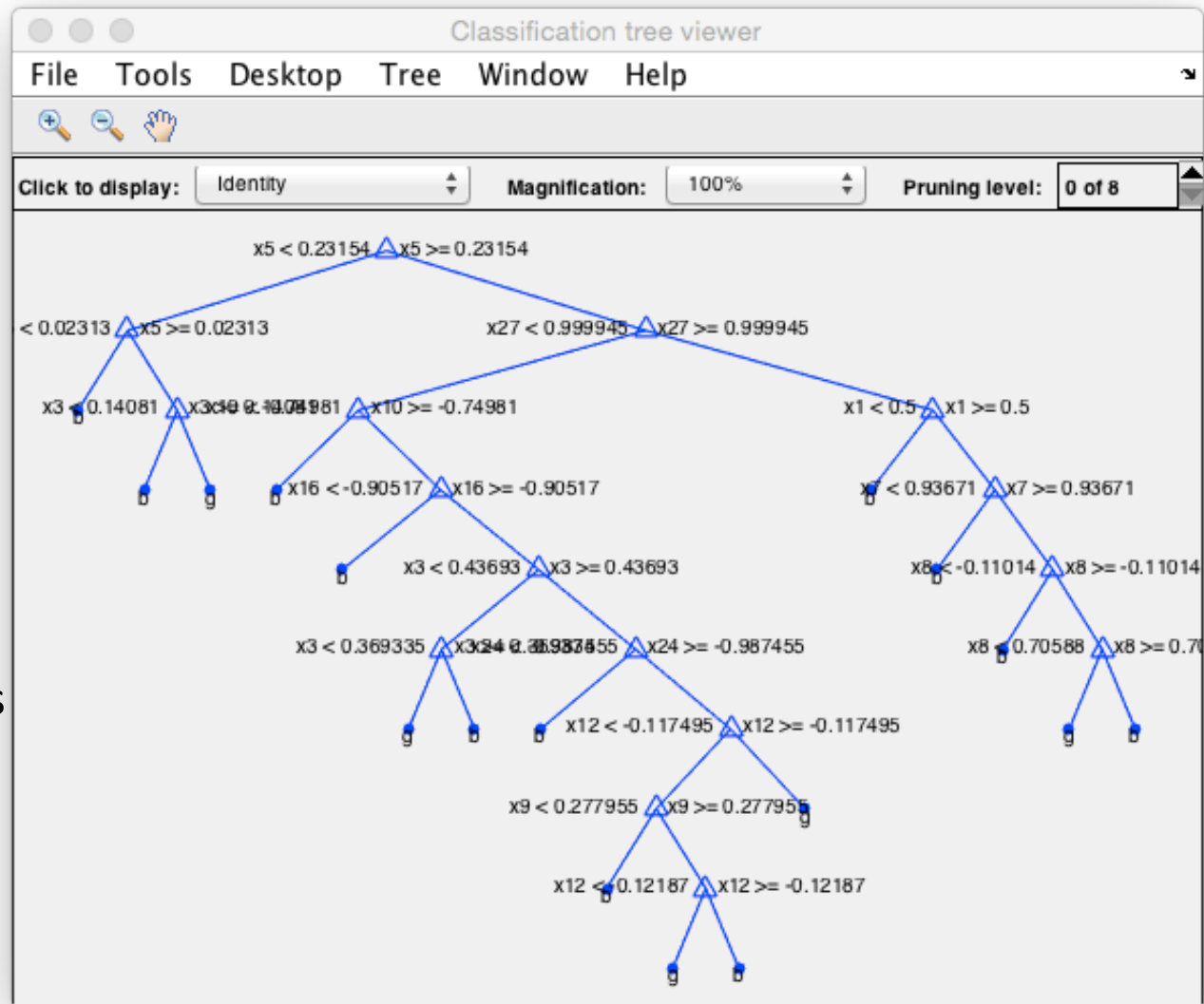
%Default MinLeafSize = 1

tc = fitctree(X,Y);

cvmodel = crossval(tc);

view(cvmodel.Trained{1},'Mode','graph')

kfoldLoss(cvmodel)



Validation error = 0.1254

Matlab example – decision tree

load ionosphere

% UCI dataset

% 34 features, 351 samples

% binary classification

rng(100)

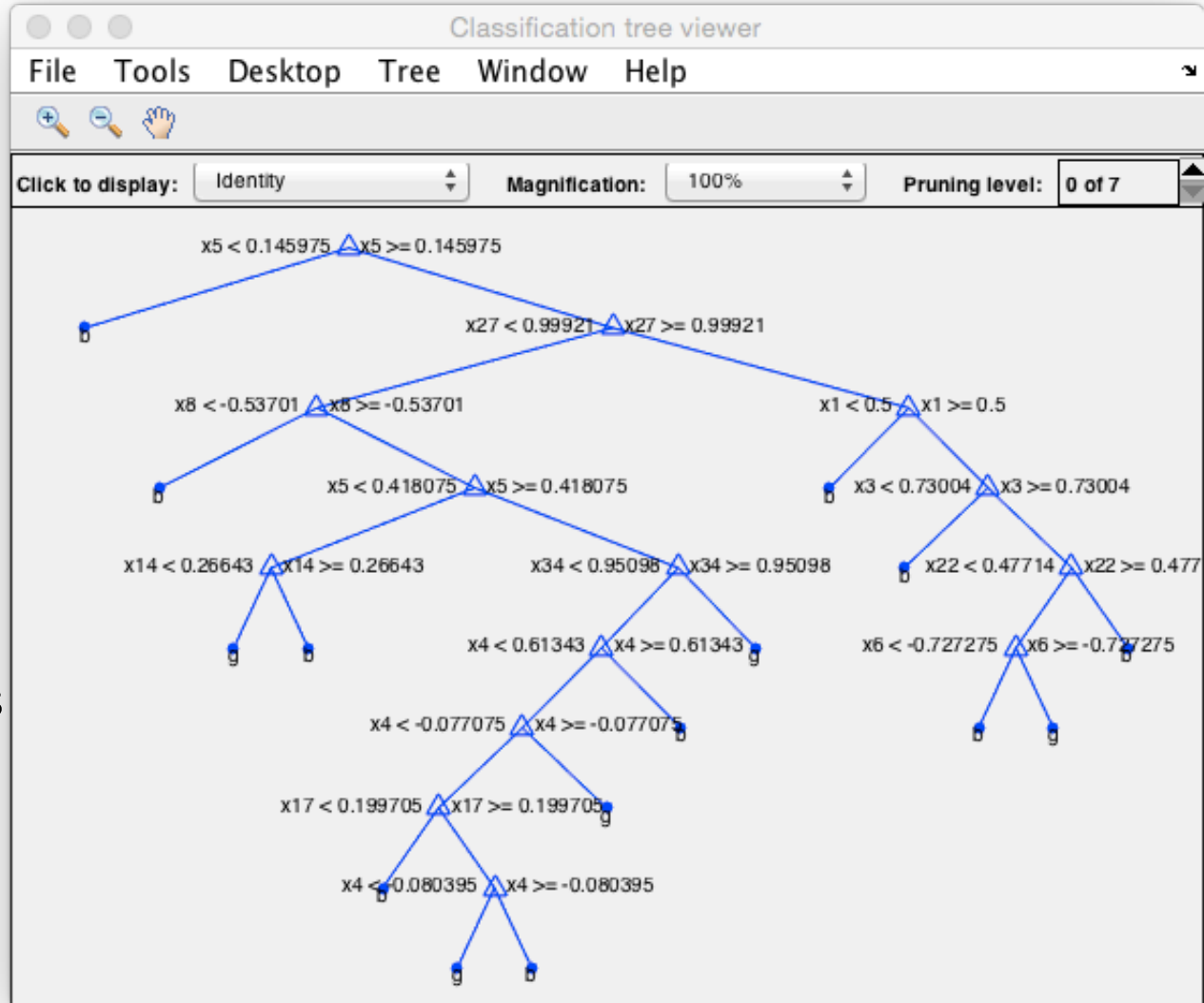
%Default MinLeafSize = 1

tc = fitctree(X,Y, 'MinLeafSize',2);

cvmodel = crossval(tc);

view(cvmodel.Trained{1},'Mode','graph')

kfoldLoss(cvmodel)



Validation error = 0.1168

Matlab example – decision tree

```
load ionosphere
```

```
% UCI dataset
```

```
% 34 features, 351 samples
```

```
% binary classification
```

```
rng(100)
```

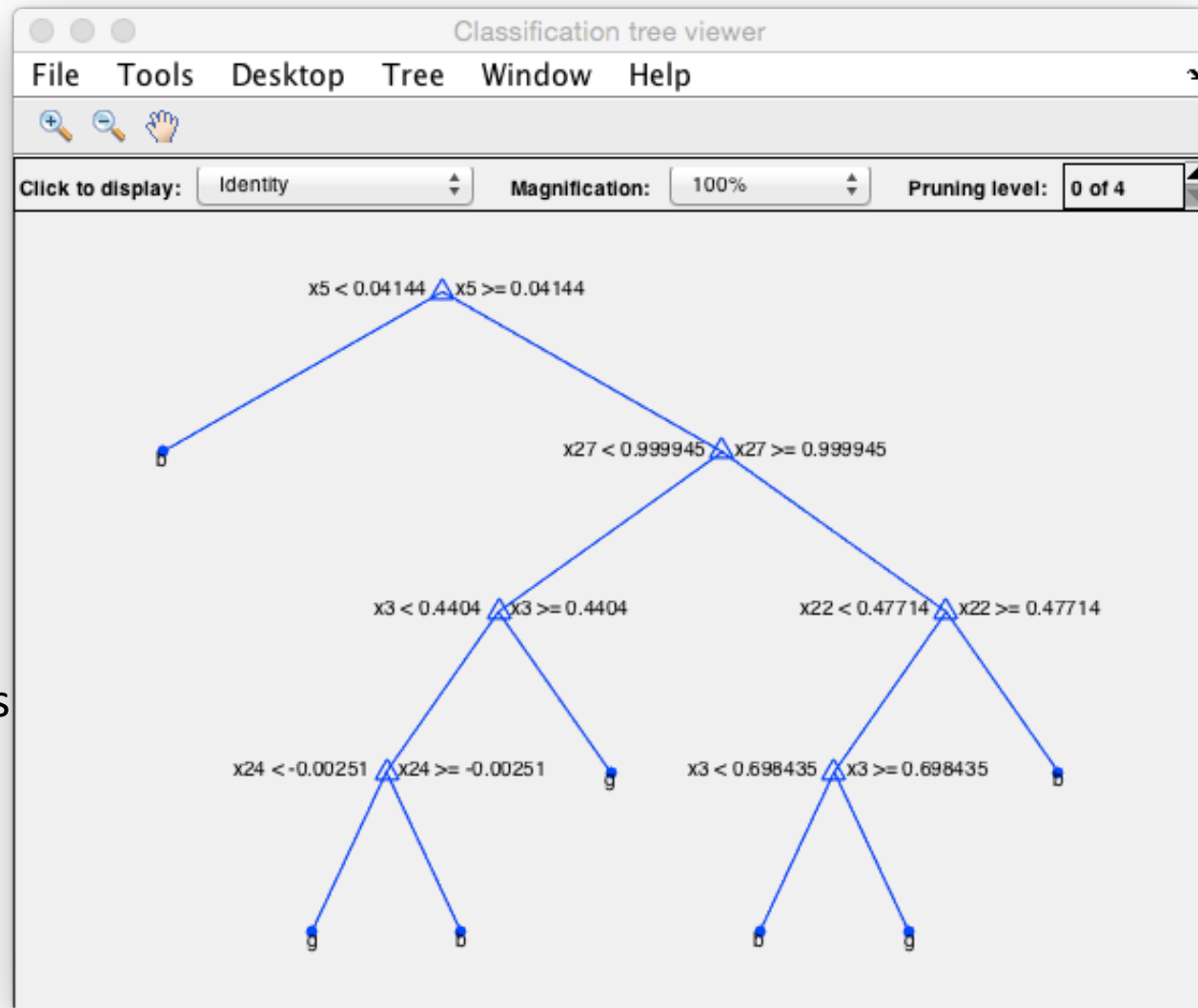
```
%Default MinLeafSize = 1
```

```
tc = fitctree(X,Y, 'MinLeafSize',10);
```

```
cvmodel = crossval(tc);
```

```
view(cvmodel.Trained{1},'Mode','graph')
```

```
kfoldLoss(cvmodel)
```



Validation error = 0.1339

Matlab example – decision trees

