# k-NN (k-Nearest Neighbor) Classifier

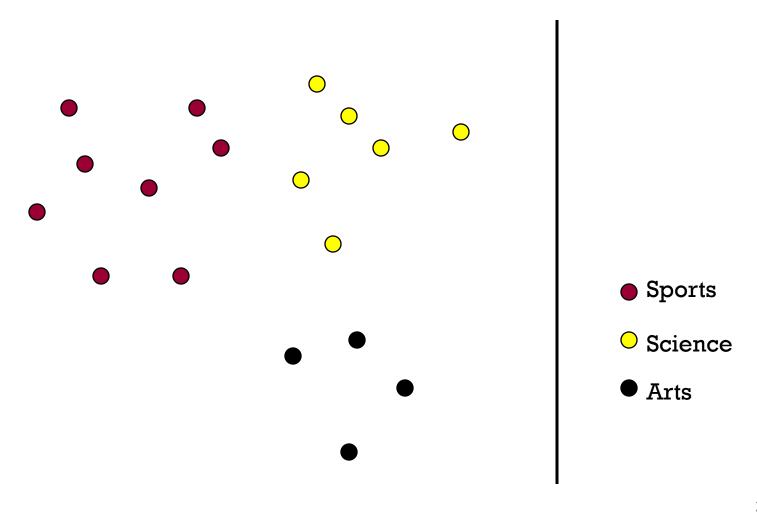
Aarti Singh

Machine Learning 10-315 Oct 14, 2019



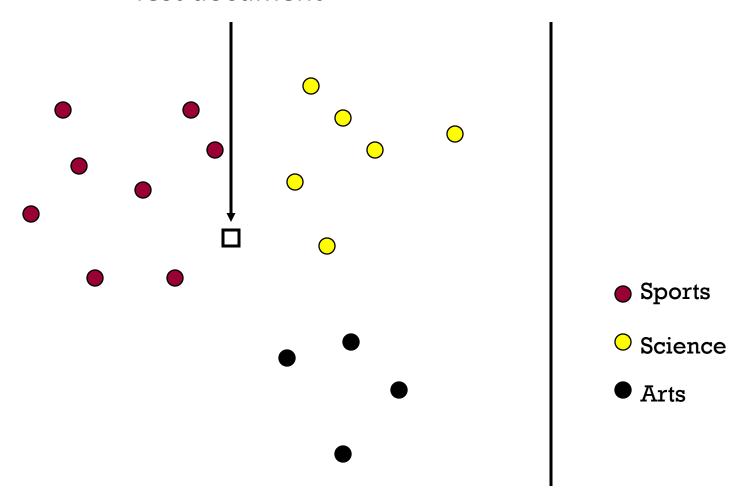


# k-NN classifier



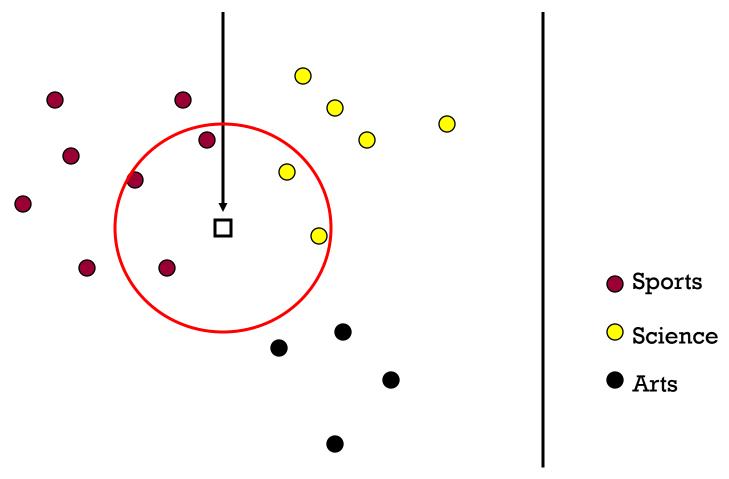
## k-NN classifier

Test document



# k-NN classifier (k=5)

Test document

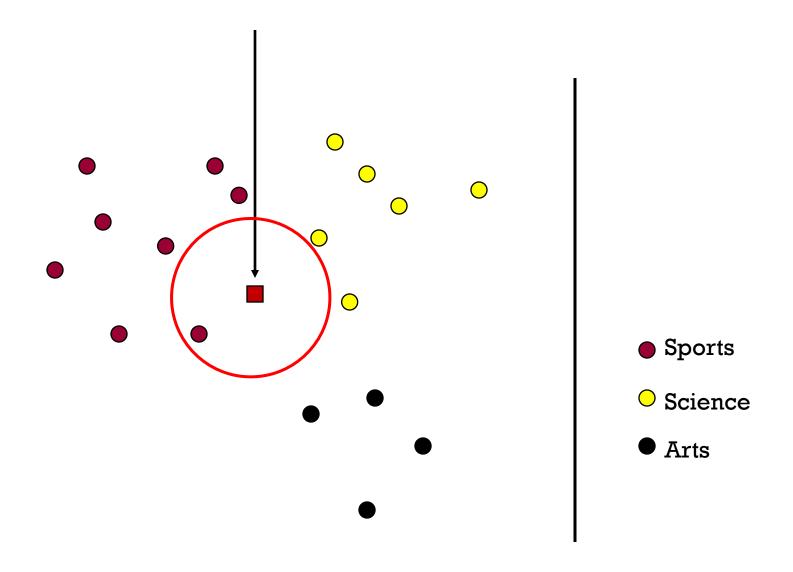


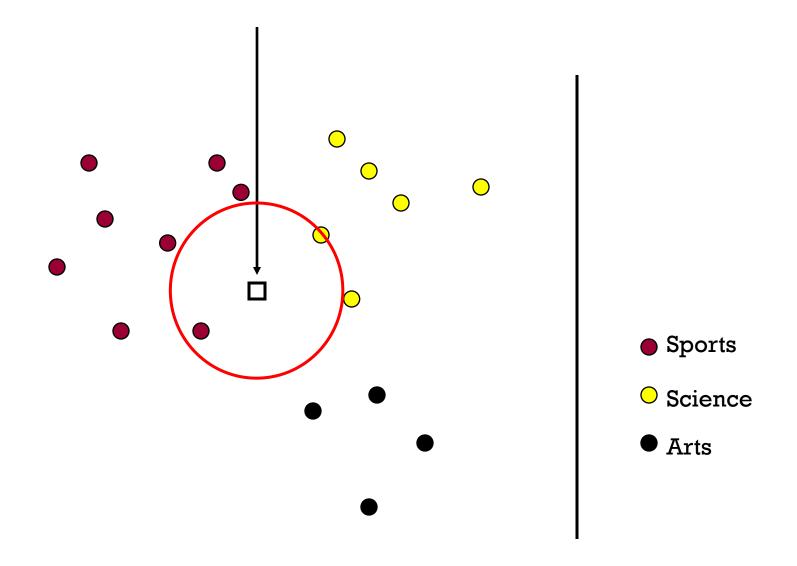
## k-NN classifier

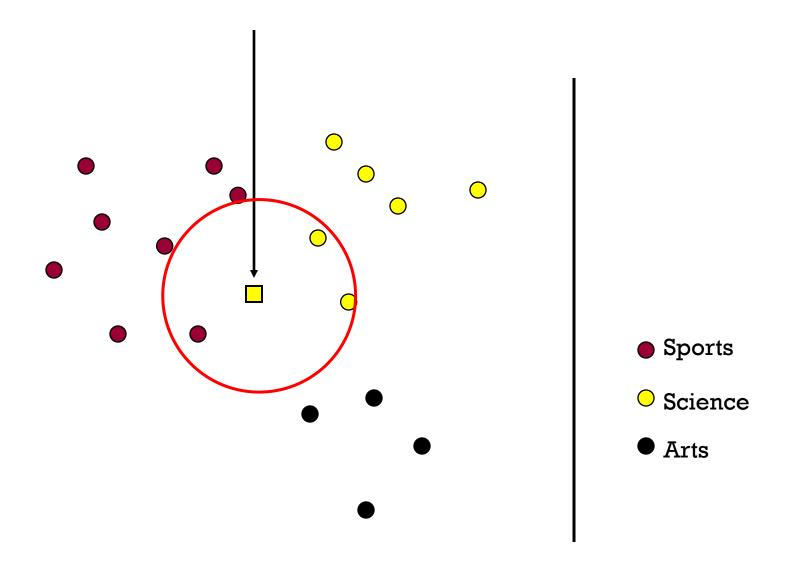
- Optimal Classifier:  $f^*(x) = \arg\max_y P(y|x)$ =  $\arg\max_y P(x|y)P(y)$
- k-NN Classifier:  $\widehat{f}_{kNN}(x) = \arg\max_{y} \ \widehat{P}_{kNN}(x|y)\widehat{P}(y)$ =  $\arg\max_{y} \ k_{y}$

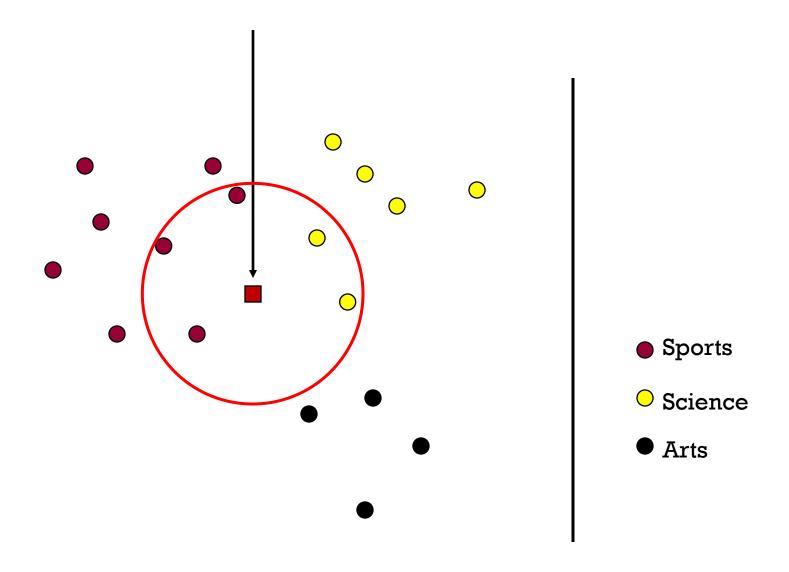
$$\widehat{P}_{kNN}(x|y) = \frac{k_y}{n_y} \longrightarrow \text{\# training pts of class y} \qquad \sum_y k_y = k$$
 amongst k NNs of x 
$$\qquad \qquad \downarrow \qquad \qquad \text{\# total training pts of class y}$$

$$\widehat{P}(y) = \frac{n_y}{n}$$



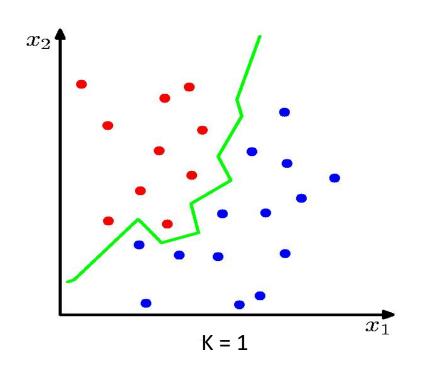




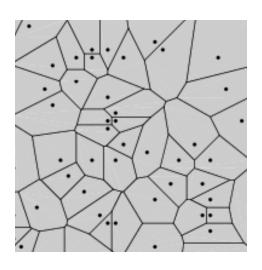


#### What is the best k?

1-NN classifier decision boundary



Voronoi Diagram



As k increases, boundary becomes smoother (less jagged).

#### What is the best k?

Approximation vs. Stability (aka Bias vs Variance) Tradeoff

- Larger K => predicted label is more stable
- Smaller K => predicted label can approximate best classifier well

# Non-parametric methods

Aka Instance-based/Memory-based learners

- Decision Trees
- k-Nearest Neighbors

#### Parametric methods

- Assume some model (Gaussian, Bernoulli, Multinomial, logistic, network of logistic units, Linear, Quadratic) with fixed number of parameters
  - Gaussian Bayes, Naïve Bayes, Logistic Regression,
     Perceptron
- Estimate parameters  $(\mu, \sigma^2, \theta, w, \beta)$  using MLE/MAP and plug in
- Pro need few data points to learn parameters
- Con Strong distributional assumptions, not satisfied in practice

#### Non-Parametric methods

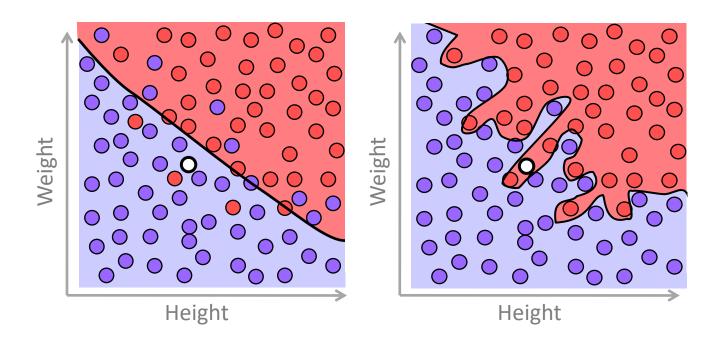
- Typically don't make any distributional assumptions
- As we have more data, we should be able to learn more complex models
- Let number of parameters scale with number of training data
- Some nonparametric methods
  - Decision Trees
  - k-NN (k-Nearest Neighbor) Classifier

# Summary

- Parametric vs Nonparametric approaches
  - Nonparametric models place very mild assumptions on the data distribution and provide good models for complex data
    - Parametric models rely on very strong (simplistic) distributional assumptions
  - Nonparametric models requires storing and computing with the entire data set.
    - Parametric models, once fitted, are much more efficient in terms of storage and computation.

# **Judging Overfitting**

# **Training Data vs. Test Data**



#### Training data

- FootballPlayer
- No
- Test data

- A good machine learning algorithm
  - Does not overfit training data
  - Generalizes well to test data

# **Training error**

Training error of a classifier f

$$\frac{1}{n} \sum_{i=1}^{n} 1_{f(X_i) \neq Y_i}$$
 Training Data 
$$\{X_i, Y_i\}_{i=1}^n$$

- What about test error?
   Can't compute it.
- How can we know classifier is not overfitting?
   Hold-out or Cross-validation

## **Hold-out method**

Can judge test error by using an independent sample of data.

#### Hold - out procedure:

n data points available  $D \equiv \{X_i, Y_i\}_{i=1}^n$ 

1) Randomly split into two sets (preserving label proportion):

Training dataset

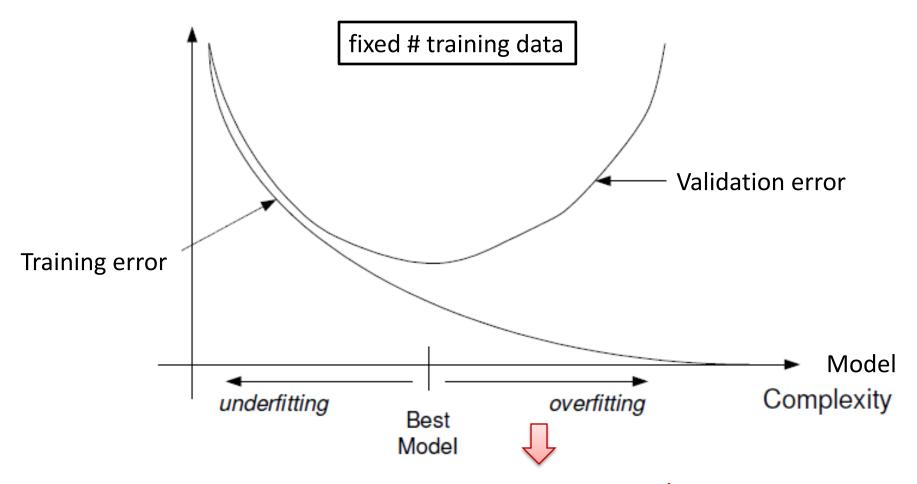
Validation/Hold-out dataset

$$D_T = \{X_i, Y_i\}_{i=1}^m \qquad D_V = \{X_i, Y_i\}_{i=m+1}^n$$

often m = n/2

2) Train classifier on  $D_T$ . Report error on validation dataset  $D_V$ . Overfitting if validation error is much larger than training error

# **Training vs. Validation Error**



Training error is no longer a good indicator of validation or test error

## **Hold-out method**

#### Drawbacks:

- May not have enough data to afford setting one subset aside for getting a sense of generalization abilities
- Validation error may be misleading (bad estimate of test error) if we get an "unfortunate" split

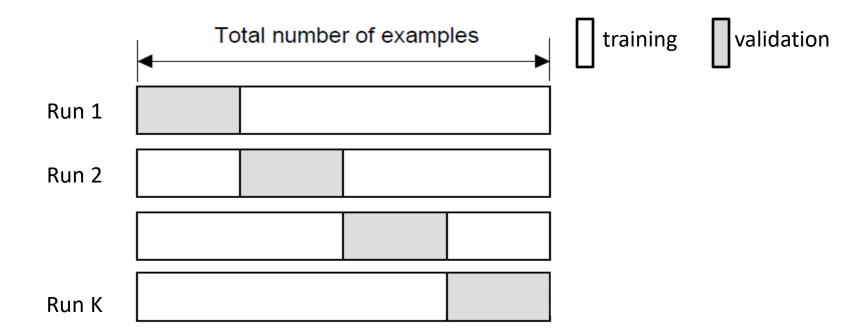
Limitations of hold-out can be overcome by a family of sub-sampling methods at the expense of more computation.

## **Cross-validation**

#### K-fold cross-validation

Create K-fold partition of the dataset.

Do K runs: train using K-1 partitions and calculate validation error on remaining partition (rotating validation partition on each run). Report average validation error



## **Cross-validation**

#### Leave-one-out (LOO) cross-validation

Special case of K-fold with K=n partitions Equivalently, train on n-1 samples and validate on only one sample per run for n runs

	4	Total number of examples	training	validation
Run 1				
Run 2				
		<b>:</b>		
Run K				

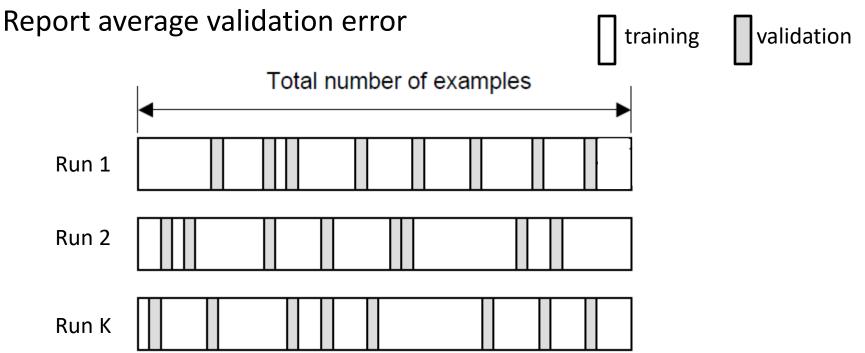
## **Cross-validation**

#### Random subsampling

Randomly subsample a fixed fraction  $\alpha n$  (0<  $\alpha$  <1) of the dataset for validation.

Compute validation error with remaining data as training data.

Repeat K times



## **Practical Issues in Cross-validation**

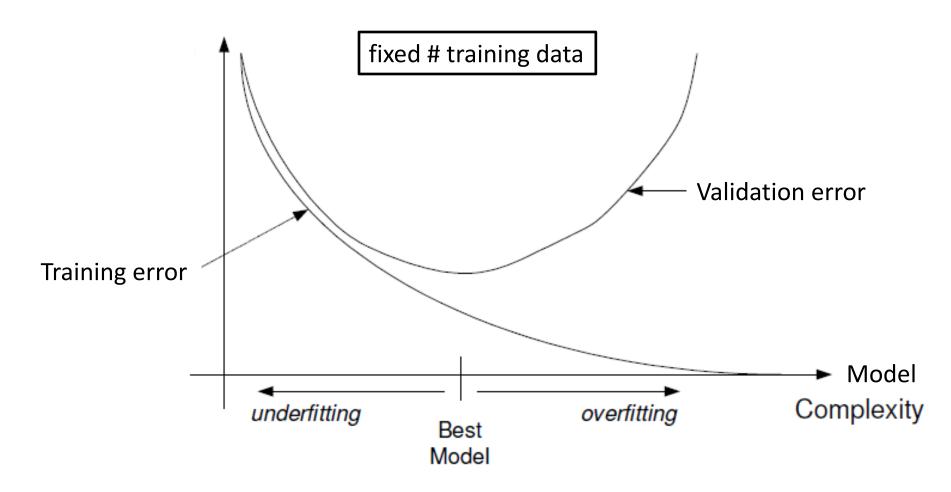
#### How to decide the values for K and a?

- Large K
  - + Validation error can approximate test error well
  - Observed validation error will be unstable (few validation pts)
  - The computational time will be very large as well (many experiments)
- Small K
  - + The # experiments and, therefore, computation time are reduced
  - + Observed validation error will be stable (many validation pts)
  - Validation error cannot approximate test error well

Common choice: K = 10,  $\alpha$  = 0.1  $\odot$ 

# **Model selection**

# **Effect of Model Complexity**



Can we select good models using hold-out or cross-validation?

# **Examples of Model Spaces**

Model Spaces with increasing complexity:

Nearest-Neighbor classifiers with increasing neighborhood sizes
 k = 1,2,3,...

Small neighborhood => Higher complexity

- Decision Trees with increasing depth k or with k leaves
   Higher depth/ More # leaves => Higher complexity
- Neural Networks with increasing layers or nodes per layer
   More layers/Nodes per layer => Higher complexity
- MAP estimates with stronger priors (larger hyper-parameters  $\beta_H$ ,  $\beta_T$  for Beta distribution or smaller variance for Gaussian prior)

How can we select the right complexity model?

# Model selection using Holdout/Cross-validation

- Train models of different complexities and evaluate their validation error using hold-out or cross-validation
- Pick model with smallest validation error (averaged over different runs for cross-validation)

Classification tree viewer Tree Window Tools Desktop Help ⊕ ⊕ 100% Click to display: Identity Magnification: Pruning level: 0 of 8  $x5 < 0.23154 \triangle x5 >= 0.23154$  $< 0.02313 \triangle x5 >= 0.02313$  $x27 < 0.999945 \triangle x27 >= 0.999945$ x3 <0.14081 x3x4e 9:49.08981 x10 >= -0.74981  $x1 < 0.5 \land x1 >= 0.5$ x16 <-0.90517 \( \text{x16} \text{ >= -0.90517} x**√** < 0.93671 Xx7 >= 0.93671 x3 < 0.43693 Xx3 >= 0.43693  $x8 < -0.11014 \times x8 > = -0.11014$ x3 < 0.369335 Xx3x24 6.385.987655 Xx24 >= -0.987455  $x8 < 0.70588 \times x8 >= 0.70$ x12 < -0.117495 x12 >= -0.117495 x9 < 0.277955 x9 >= 0.27795 $x12 < 0.12187 \times 12 >= -0.12187$ 

load ionosphere
% UCI dataset
% 34 features, 351 samples
% binary classification
rng(100)

%Defaulty MinLeafSize = 1
tc = fitctree(X,Y);
cvmodel = crossval(tc);
view(cvmodel.Trained{1},'Mode','graph')
kfoldLoss(cvmodel)

Validation error = 0.1254

⊕ ⊕ ₩ Identity 100% Click to display: Magnification: Pruning level: 0 of 7  $x5 < 0.145975 \triangle x5 >= 0.145975$  $x27 < 0.99921 \Delta x27 >= 0.99921$ x8 <-0.53701 Ax8 >= -0.53701  $x1 < 0.5 \land x1 >= 0.5$  $x5 < 0.418075 \Delta x5 >= 0.418075$ x3 < 0.73004 \( \text{\tin}\text{\tinx}\text{\te}\tint{\texi}\text{\text{\text{\text{\text{\text{\text{\text{\texi}\text{\text{\texit{\text{\texi}\text{\text{\text{\text{\texi}\text{\texit{\ti}\xint{\text{\texit{\text{\texi}\text{\text{\texit{\text{\text{ x14 < 0.26643 Ax14 >= 0.26643 x22 < 0.47714 x22 >= 0.477 x34 < 0.95098 Ax34 >= 0.95098 x4 < 0.61343 X4 >= 0.61343 x6 < -0.727275 //x6 >= -0.7 27275 x4 < -0.077075 Xx4 >= -0.077075 x17 < 0.199705 Xx17 >= 0.199705 x4 40.080395 x4 >= -0.080395 %Defaulty MinLeafSize = 1

Classification tree viewer

Help

Window

```
load ionosphere
% UCI dataset
% 34 features, 351 samples
% binary classification
rng(100)
```

tc = fitctree(X,Y, 'MinLeafSize',2); cvmodel = crossval(tc); view(cvmodel.Trained{1},'Mode','graph') kfoldLoss(cvmodel)

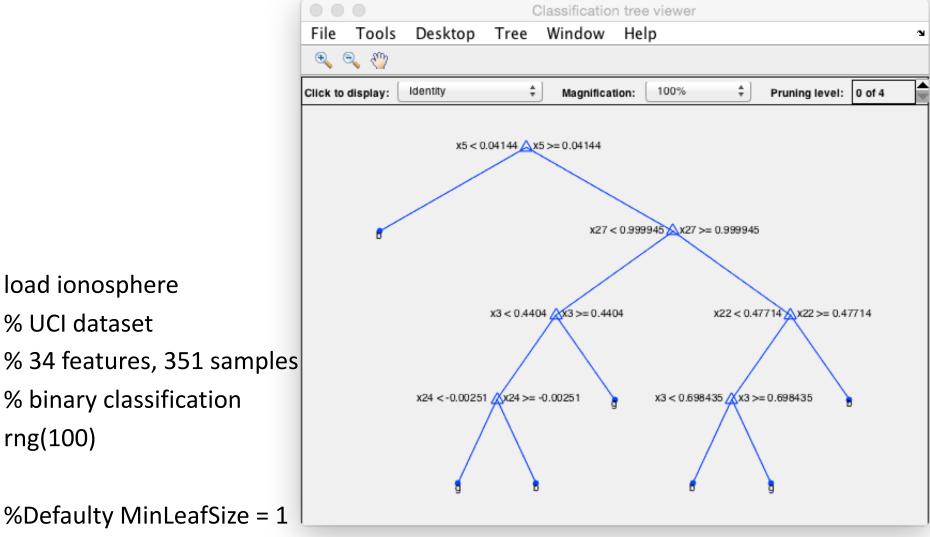
File

Tools

Desktop

Tree

Validation error = 0.1168



load ionosphere % UCI dataset % 34 features, 351 samples % binary classification rng(100)

tc = fitctree(X,Y, 'MinLeafSize',10); cvmodel = crossval(tc); view(cvmodel.Trained{1},'Mode','graph') kfoldLoss(cvmodel)

Validation error = 0.1339

