Machine Learning and Data Mining

Nearest neighbor methods

Prof. Alexander Ihler







Supervised learning

V

- **Notation**
 - Features X
 - Targets
 - Predictions \hat{y}
 - Parameters θ



Regression; Scatter plots



- Suggests a relationship between x and y
- Regression: given new observed x^(new), estimate y^(new)

Nearest neighbor regression



"Predictor": Given new features: Find nearest example Return its value

• Find training datum $x^{(i)}$ closest to $x^{(new)}$; predict $y^{(i)}$

Nearest neighbor regression



"Predictor": Given new features: Find nearest example Return its value

- Find training datum x⁽ⁱ⁾ closest to x^(new); predict y⁽ⁱ⁾
- Defines an (implict) function f(x)
- "Form" is piecewise constant









More Data Points



 $X_1 \rightarrow$

More Complex Decision Boundary



Machine Learning and Data Mining

Nearest neighbor methods: K-Nearest Neighbors

Prof. Alexander Ihler



+





K-Nearest Neighbor (kNN) Classifier

- Find the k-nearest neighbors to <u>x</u> in the data
 - i.e., rank the feature vectors according to Euclidean distance
 - select the k vectors which are have smallest distance to \underline{x}
- Regression
 - Usually just average the y-values of the k closest training examples
- Classification
 - ranking yields k feature vectors and a set of k class labels
 - pick the class label which is most common in this set ("vote")
 - classify <u>x</u> as belonging to this class
 - Note: for two-class problems, if k is odd (k=1, 3, 5, ...) there will never be any "ties"; otherwise, just use (any) tie-breaking rule
- "Like" the optimal estimator, but using nearest k points to estimate p(y|x)
- "Training" is trivial: just use training data as a lookup table, and search to classify a new datum

kNN Decision Boundary

- Piecewise linear decision boundary
- Increasing k "simplifies" decision boundary
 - Majority voting means less emphasis on individual points







kNN Decision Boundary

- Piecewise linear decision boundary
- Increasing k "simplifies" decision boundary
 - Majority voting means less emphasis on individual points







kNN Decision Boundary

- Piecewise linear decision boundary
- Increasing k "simplifies" decision boundary
 - Majority voting means less emphasis on individual points



K = 25



Complexity & Overfitting

- Complex model predicts all training points well
- Doesn't generalize to new data points
- k = 1 : perfect memorization of examples (complex)
- k = m : always predict majority class in dataset (simple)
- Can select k using validation data, etc.



K-Nearest Neighbor (kNN) Classifier

- Theoretical Considerations
 - as k increases
 - we are averaging over more neighbors
 - the effective decision boundary is more "smooth"
 - as N increases, the optimal k value tends to increase
 - k=1, m increasing to infinity : error < 2x optimal
- Extensions of the Nearest Neighbor classifier
 - Weighted distances $d(x, x') = \sqrt{\sum_i w_i (x_i x'_i)^2}$
 - e.g., some features may be more important; others may be irrelevant
 - Mahalanobis distance: $d(x, x') = \sqrt{(x x') \cdot S^{-1} \cdot (x x')}$
 - Fast search techniques (indexing) to find k-nearest points in d-space
 - Weighted average / voting based on distance

Summary

- K-nearest neighbor models
 - Classification (vote)
 - Regression (average or weighted average)
- Piecewise linear decision boundary
 - How to calculate
- Test data and overfitting
 - Model "complexity" for knn
 - Use validation data to estimate test error rates & select k