k-Nearest Neighbors

Samuel Taylor









Chervny

Modry

Chervny

Modry

Chervny



Outline

- Introduction
- Algorithm
- How good is it?
- Application



Are you familiar with...

- Approximation–generalization tradeoff?
- VC dimension?
- Regularization?

Outline

- Introduction
- Algorithm
- How good is it?
- Application

How did we do that?



- Observe training data
- Find most similar datum
 - Color
 - Points
- Pick the class of the similar datum



Number of points

The Nearest Neighbor Algorithm

- Given a training set $D = \{(x_1, y_1), \dots, (x_N, y_N)\}$ and a data point u,
- Find the point $d = (x_i, y_i)$ in **D** for which x_i is closest to u
- Return y_j

The k-Nearest Neighbor Algorithm

- Given a training set $D = \{(x_1, y_1), \dots, (x_N, y_N)\}$ and a data point u,
- Find the *k* nearest points to $u(x_{n1}, y_{n1}), \dots, (x_{nk}, y_{nk})$ in **D**
- Return sign(sum(y_{ni} for *i* in range(k)) / k)

Observations

- No training required
- Low values of k tend toward overfit, high values tend to underfit







Observations

- Very low in-sample error
- Infinite VC-dimension

Outline

- Introduction
- Algorithm
- How good is it?
- Application

Pretty good

- $\leq 2 \times$ minimum possible out of sample error
 - Probably, as N -> ∞
 - Under reasonable assumptions

Proving $E_{out} \leq 2E_{out}^*$ $\pi(\mathbf{x}) = \mathbb{P}[y = +1 | \mathbf{x}].$ \leftarrow the target in logistic regression Assume $\pi(\mathbf{x})$ is continuous and $\mathbf{x}_{[1]} \xrightarrow{N \to \infty} \mathbf{x}$. Then $\pi(\mathbf{x}_{[1]}) \xrightarrow{N \to \infty} \pi(\mathbf{x})$. $\mathbb{P}[q_N(\mathbf{x}) \neq y] = \mathbb{P}[y = +1, y_{[1]} = -1] + \mathbb{P}[y = -1, y_{[1]} = +1],$ $= \pi(\mathbf{x}) \cdot (1 - \pi(\mathbf{x}_{[1]})) + (1 - \pi(\mathbf{x})) \cdot \pi(\mathbf{x}_{[1]}),$ $\rightarrow \pi(\mathbf{x}) \cdot (1 - \pi(\mathbf{x})) + (1 - \pi(\mathbf{x})) \cdot \pi(\mathbf{x}),$

 $= 2\pi(\mathbf{x}) \cdot (1 - \pi(\mathbf{x})),$

 $\leq 2\min\{\pi(\mathbf{x}), 1-\pi(\mathbf{x})\}.$

The best you can do is

 $E_{\text{out}}^*(\mathbf{x}) = \min\{\pi(\mathbf{x}), 1 - \pi(\mathbf{x})\}.$

© 🎢 Creator: Malik Magdon-Ismail

Similarity and Nearest Neighbor: 9 /16

NN is self-regularizing











Outline

- Introduction
- Algorithm
- How good is it?
- Application



url = 'http://dash.samueltaylor.org/commute_time.csv'
commutes = pd.read_csv(url, index_col=0, parse_dates=True)
X = 60 * (commutes.index.hour * 60 + commutes.index.minute)
y = commutes['duration']

X_train, X_test, y_train, y_test = train_test_split(X, y, test size=0.33)

regr = neighbors.KNeighborsRegressor(3)
regr.fit(X train, y train)

predictions = regr.predict(X_test)



Choosing k

- Controls the tradeoff between approximation and generalization
- 3
- floor(sqrt(N))
 - As $N \to \infty$ and $k(N)/N \to 0$, out of sample error approaches optimal





Choosing k

- Controls the tradeoff between approximation and generalization
- 3
- floor(sqrt(N))
 - As $N \to \infty$ and $k(N)/N \to 0$, out of sample error approaches optimal
- Cross validation

Unaddressed questions

- How do I choose a similarity function?
- How do I choose features?
- Why isn't the scikit implementation terribly slow (esp. on large datasets)?

Further reading

- Learning From Data; Abu-Mostafa, Magdon-Ismail, and Lin
- Similarity and Nearest Neighbor; Malik Magdon-Ismail



Samuel Taylor

kdd@samueltaylor.org