

Day 15 - 1 November - Cross Validation and K-Nearest Neighbor

Agenda:

- Cross Validation
- K-Nearest Neighbor

Model Validation

Suppose you have multiple ^{trained} predictors you are choosing between. How do you select the best one?

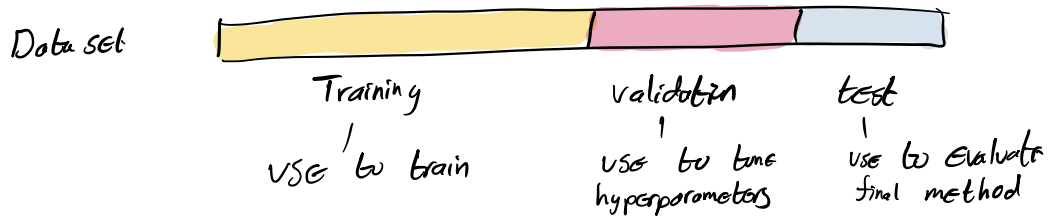
eg Ridge regression

$$\min_{\theta} \|y - X\theta\|^2 + \lambda \|\theta\|^2$$

parameter

Which value of λ should you choose?
hyperparameter

Ideally, use validation data



$$\min_{\lambda} \frac{1}{m} \sum_{i=1}^m \mathcal{L}(y_i^{\text{val}}, f_{\lambda}(x_i^{\text{val}}))$$

predictor w/
hyperparameter λ

Challenges:

- Need data for validation,
- Need independent data for test
- data is often expensive

K-fold Cross Validation

Chap 7 of Hastie

7.10.1 K-Fold Cross-Validation

Ideally, if we had enough data, we would set aside a validation set and use it to assess the performance of our prediction model. Since data are often scarce, this is usually not possible. To finesse the problem, **K-fold cross-validation uses part of the available data to fit the model, and a different part to test it. We split the data into K roughly equal-sized parts;** for example, when $K = 5$, the scenario looks like this:

/"folds"

242 7. Model Assessment and Selection

1	2	3	4	5
Train	Train	Validation	Train	Train

For the k th part (third above), we fit the model to the other $K - 1$ parts of the data, and calculate the prediction error of the fitted model when predicting the k th part of the data. We do this for $k = 1, 2, \dots, K$ and combine the K estimates of prediction error.

Here are more details. Let $\kappa : \{1, \dots, N\} \mapsto \{1, \dots, K\}$ be an indexing function that indicates the partition to which observation i is allocated by the randomization. Denote by $\hat{f}^{-\kappa}(x)$ the fitted function, computed with the k th part of the data removed. Then the cross-validation estimate of prediction error is

$$CV(\hat{f}) = \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{f}^{-\kappa(i)}(x_i)). \quad (7.48)$$

avg over
each
data point

train on all
data not in
the fold containing
 i^{th} data point

Equivalently: average over K folds of the average loss on each fold when trained on all other folds

Hyper
parameter
tuning

Typical choices of K are 5 or 10 (see below). The case $K = N$ is known as *leave-one-out* cross-validation. In this case $\kappa(i) = i$, and for the i th observation the fit is computed using all the data except the i th.

Given a set of models $f(x, \alpha)$ indexed by a tuning parameter α , denote by $\hat{f}^{-k}(x, \alpha)$ the α th model fit with the k th part of the data removed. Then for this set of models we define

$$\text{CV}(\hat{f}, \alpha) = \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{f}^{-\kappa(i)}(x_i, \alpha)). \quad (7.49)$$

The function $\text{CV}(\hat{f}, \alpha)$ provides an estimate of the test error curve, and we find the tuning parameter $\hat{\alpha}$ that minimizes it. Our final chosen model is $f(x, \hat{\alpha})$, which we then fit to all the data.

How to choose # of folds?

If $K=N$, • have to solve N problems
⇒ EXPENSIVE.
• low bias, high variance

If K too small, not enough data used
for training model.

Compromise $k = 5$ or 10

7.10.2 *The Wrong and Right Way to Do Cross-validation*

Consider a classification problem with a large number of predictors, as may arise, for example, in genomic or proteomic applications. A typical strategy for analysis might be as follows:

1. Screen the predictors: find a subset of “good” predictors that show fairly strong (univariate) correlation with the class labels
2. Using just this subset of predictors, build a multivariate classifier.
3. Use cross-validation to estimate the unknown tuning parameters and to estimate the prediction error of the final model.

Is this a correct application of cross-validation?

The predictor screening step uses all of the data. When training each model for cross-validation, this means each of those models will have trained on the test data that was left out.

Here is the correct way to carry out cross-validation in this example:

1. Divide the samples into K cross-validation folds (groups) at random.
2. For each fold $k = 1, 2, \dots, K$
 - (a) Find a subset of “good” predictors that show fairly strong (univariate) correlation with the class labels, using all of the samples except those in fold k .
 - (b) Using just this subset of predictors, build a multivariate classifier, using all of the samples except those in fold k .
 - (c) Use the classifier to predict the class labels for the samples in fold k .

The error estimates from step 2(c) are then accumulated over all K folds, to produce the cross-validation estimate of prediction error. The lower panel

In general, with a multistep modeling procedure, cross-validation must be applied to the entire sequence of modeling steps. In particular, samples must be “left out” before any selection or filtering steps are applied. There is one qualification: initial *unsupervised* screening steps can be done before samples are left out. For example, we could select the 1000 predictors with highest variance across all 50 samples, before starting cross-validation. Since this filtering does not involve the class labels, it does not give the predictors an unfair advantage.

K nearest neighbors (KNN)

- Method for regression & classification
- Idea: Find k closest training examples and average them / do a majority vote

Classification

Training data: $\{(X_i, y_i)\}$ w/ $X_i \in \mathbb{R}^d$, $y_i \in \{1 \dots M\}$

Predictor:

$$P(y=c | \underset{\substack{\text{train data} \\ S}}{X}, k) = \frac{1}{k} \sum_{i \in N_k} \mathbb{1}_{y_i=c}$$

S

$N_k =$ set of indices i
of the k nearest
 X_i to X

$\hat{y}(X) =$ most common value
of y_i among k nearest
values of X_i to X .

Distance measured by l_2 norm (or
alternative of your choice)

2 class
example:

$k=1$

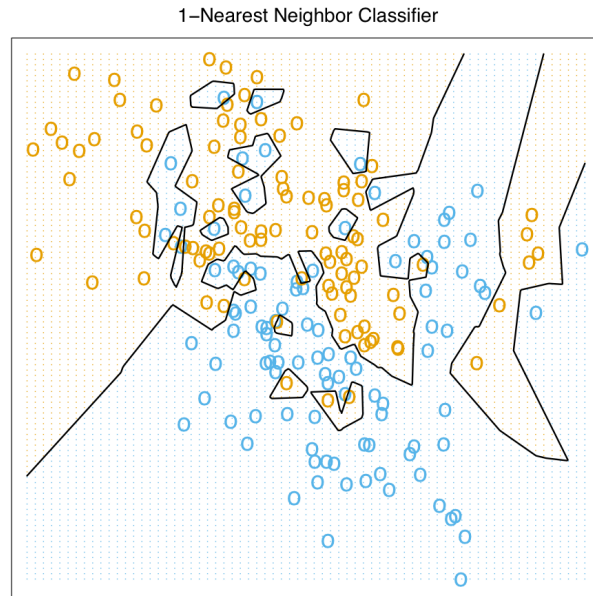
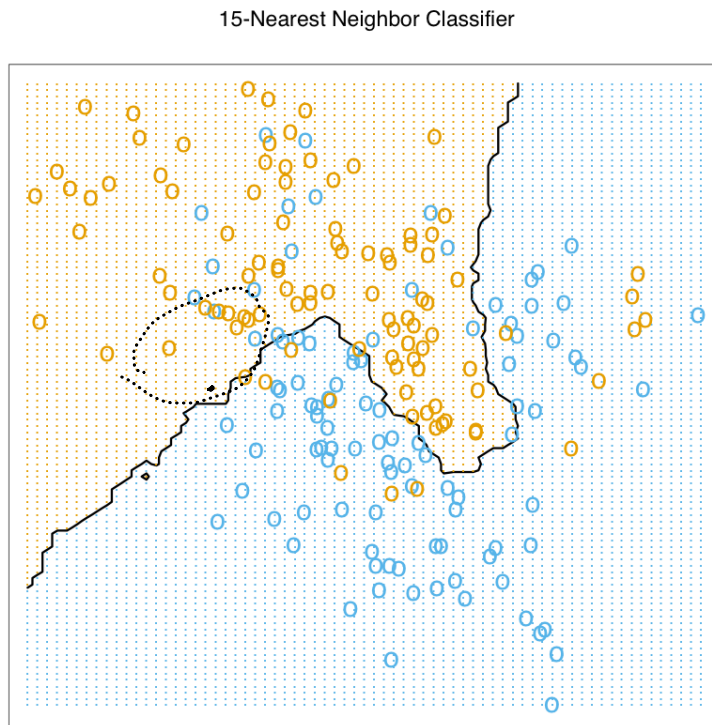


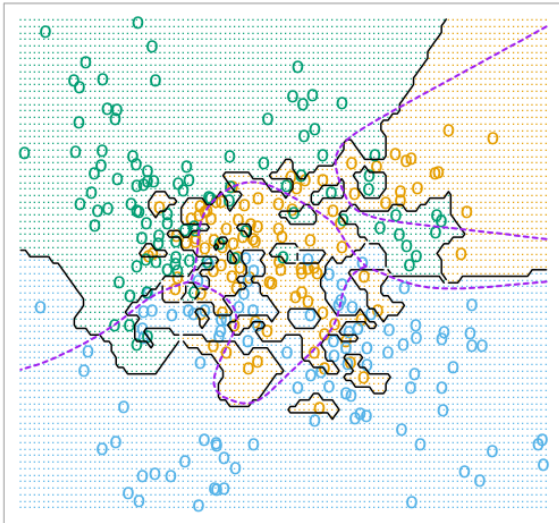
FIGURE 2.3. The same classification example in two dimensions as in Figure 2.1. The classes are coded as a binary variable (BLUE = 0, ORANGE = 1), and then predicted by 1-nearest-neighbor classification.

$k=15$

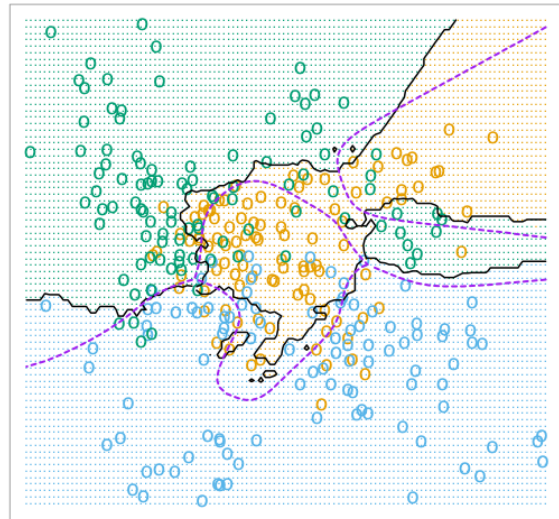


3 class example

1-Nearest Neighbor



15-Nearest Neighbors



— Decision boundary for Bayes Classifier

Bayes Classifier

Assume (X, y) follows a joint distribution

$$\hat{y}(x) = \underset{y}{\operatorname{argmax}} P(y | x)$$

/
Bayes Classifier gives most likely class of y given x . Can't compute this as underlying dist is unknown.

KNN and bias-variance tradeoff

What is the effect of different values of k ?

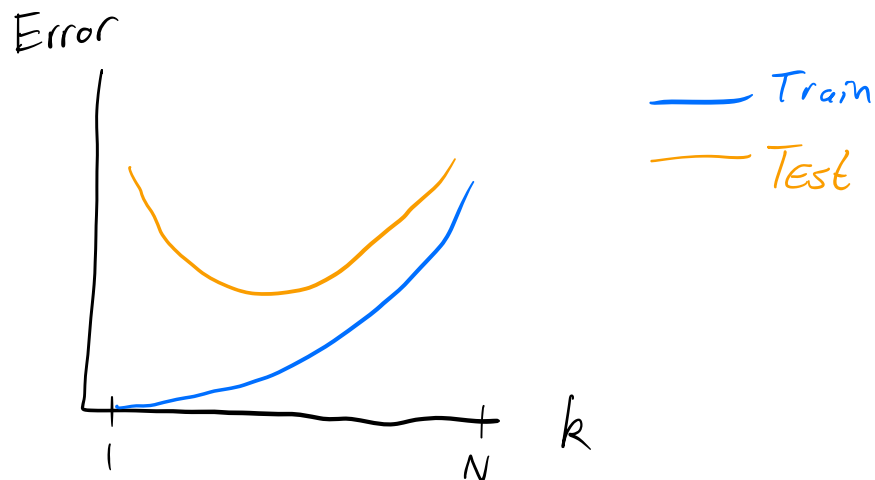
If k too small:

Bias: High or low
Variance: High or low

If k too large?

Bias: High or low
Variance: High or low

Plot: Train Error vs k & Test Error vs k



How do you find k ?

Regression with k nearest neighbors

Training data: $\{(X_i, y_i)\}_{i=1, \dots, n}$ w/ $X_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}^m$

Prediction: $\hat{y} = \frac{1}{k} \sum_{i \in N_k} y_i$

N_k = set of indices i
of the k nearest
 X_i to x

Bias-Variance tradeoff for regression:

Consider a model $y_i = h(x_i) + \varepsilon_i$ w/ $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$
|
true response function

$$\begin{aligned} \text{Bias: } \left(\mathbb{E}_y \hat{y}(x) - h(x) \right)^2 &= \left(\frac{1}{k} \sum_{i \in N_k} \mathbb{E}_y(y_i) - h(x) \right)^2 \\ &= \left(\frac{1}{k} \sum_{i \in N_k} h(x_i) - h(x) \right)^2 \end{aligned}$$

If k fixed and $n \rightarrow \infty$,

$$\frac{1}{k} \sum_{i \in N_k} h(x_i) \approx \underbrace{h(x)}_{\text{low bias}}$$

If $k = n$ and $n \rightarrow \infty$ high bias

$$\frac{1}{k} \sum_{i \in N_k} h(x_i) \approx \text{average of } h$$

Variance^o

$$\text{Var}(\hat{y}(x)) = \text{Var}\left[\frac{1}{k} \sum_{i \in N_k} y_i\right]$$

$$= \frac{1}{k^2} \text{Var} \sum_{i \in N_k} y_i$$

$$= \frac{1}{k^2} k \text{Var}(y_i)$$

$$= \frac{\sigma^2}{k}$$

Small k ^o high variance

Large k ^o low variance

ISSUE: KNN may fail w/ high dimensional data.

Why? Because "all points in high dimensional space are roughly equidistant"

CURSE OF DIMENSIONALITY

Illustration:

Consider a sphere of radius 1 in \mathbb{R}^d . What fraction of the sphere's volume is located inside radius $1-\epsilon$?

$$\frac{(1-\epsilon)^d}{1^d} = (1-\epsilon)^d$$

which is very small for large d !

A randomly generated point in a ball of radius r in high dim space is almost certainly near its boundary.