

# Project 2: Predictive Modelling

STAT GU4243

*Applied Data Science*

Cynthia Rush  
Columbia University

February 19, 2018

# SUPERVISED LEARNING

Data:

$$\begin{aligned} &(y_1, x_{11}, x_{12}, \dots, x_{1p}) \\ &(y_2, x_{21}, x_{22}, \dots, x_{2p}) \\ &\dots \\ &(y_n, x_{n1}, x_{n2}, \dots, x_{np}) \end{aligned}$$

$Y$  referred to as ‘**outputs**’, ‘**responses**’, or ‘**dependent variables**’

$X$  referred to as ‘**inputs**’, ‘**predictors**’, ‘**features**’, or ‘**independent variables**’

# SUPERVISED LEARNING

In supervised learning, outcome variable  $Y$  is given.

In unsupervised learning, there is no label  $Y$ .

## Learning Tasks

1. **Prediction:** Use  $X$  to construct a model to predict  $Y$ .
2. **Inference:** Identify subject-matter knowledge by understanding the learned model.

For categorical  $Y$ , it is referred to as a ‘classification’ problem.

For quantitative  $Y$ , often called ‘regression’.

# ASSUMPTIONS AND TERMINOLOGY

In a **classification problem**, we record measurements

$$X_1 = (x_{11}, x_{12}, \dots, x_{1p}), X_2 = (x_{21}, x_{22}, \dots, x_{2p}), \dots$$

We assume:

1. All measurements can be represented as elements of  $\mathbb{R}^p$  ( $p$  dimensional Euclidean space).
2. Each  $x_i$  belongs to exactly one out of  $K$  categories, called **classes**. We express this using variables  $y_i \in [K]$ , called **class labels**:

$$y_i = k \quad \leftrightarrow \quad "x_i \text{ in class } k"$$

3. The classes are characterized by the (unknown!) joint distribution of  $(X, Y)$ , whose density we denote  $p(x, y)$ .
4. The only information available on the distribution  $p$  is a set of example measurements *with* labels,

$$(y_1, x_{11}, x_{12}, \dots, x_{1p}), \dots, (y_n, x_{n1}, x_{n2}, \dots, x_{np}),$$

called the **training data**.

## Definition

A **classifier** is a function

$$f : \mathbb{R}^p \rightarrow [K] ,$$

i.e. a function whose argument is a measurement and whose output is a class label (one of  $1, 2, \dots, K$ ).

## Learning task

Using the training data, we have to estimate a good classifier. This estimation procedure is also called **training**.

A good classifier should generalize well to new data. Ideally, we would like it to perform with high accuracy on data sampled in the same way as the training data (i.e. also from  $p(x, y)$ ).

## Simplifying assumption

We can consider the two-class case ( $K=2$ ), which is also called **binary classification**. In this case, we use the notation

$$Y \in \{-1, +1\} \quad \text{instead of} \quad Y \in \{1, 2\}$$

# CLASSIFICATION

Most classification methods use:

1. Linear or non-linear **decision boundaries** between classes.
2. **Discriminant functions**: for each class  $k \in [K]$ , define  $\hat{f}_k(x)$ .
  - ▶ Prediction may then be computed as:  $\arg \max_k \hat{f}_k(x)$ .
  - ▶ Most current methods work in this domain.

## What We Won't Talk About

Can't cover every possible classification method you may need to use for the project – there are too many and I don't know them all

## What We Will Talk About

- ▶ Framework for comparing classification methods: how do we judge performance?
  - ▶ Loss functions and risk
  - ▶ Test error vs. training error
  - ▶ Cross-validation and bootstrap
- ▶ General properties of classification methods: bias vs variance, complexity, curse of dimensionality, etc.
- ▶ Some classification basics: logistic regression, linear discriminant analysis, kNN
- ▶ Some strategies for building better classifiers: boosting, bagging, etc.



# Some Lite Decision Theory

# LOSS FUNCTIONS

First consider:

- ▶ Real-valued, random input  $X \in \mathbb{R}^p$
- ▶ Real-valued, random  $Y \in \mathbb{R}$  (quantitative, not categorical  $Y$  for now).

We assume a joint probability distribution  $p(x, y)$ .

## Definition

Want a function  $f(X)$  for predicting  $Y$  given values of  $X$ . This theory requires a loss function:

$$L(Y, f(X)) : \mathbb{R} \times \mathbb{R} \rightarrow [0, \infty)$$

for penalizing prediction errors.

## Squared Error

The most common and convenient is squared error loss:

$$L(Y, f(X)) = (Y - f(X))^2.$$

# EXPECTED PREDICTION ERROR

## Motivation

It may be a good strategy to allow (even expensive) errors for values of  $X$  which are very unlikely to occur.

## Definition

The expected prediction error,  $\text{EPE}(f)$ , of a classifier  $f$  is its expected loss under  $p$ , that is,

$$\text{EPE}(f) = \mathbb{E}_{(X,Y)}[L(Y, f(X))].$$

# LOSS FUNCTIONS

## Example

For squared error loss,

$$\text{EPE}(f) = \mathbb{E}_{(X,Y)}(Y - f(X))^2 = \mathbb{E}_X \mathbb{E}_{Y|X}((Y - f(X))^2 | X)$$

Using the above,  $f(x) = \mathbb{E}(Y|X = x)$  minimizes  $\text{EPE}(f)$ , which is known as the regression function.

## Interpretation

Tells us, best prediction of  $Y$  when  $X = x$  is given by  $\mathbb{E}(Y|X = x)$ , when best is measured by square error loss.

Therefore, the type of loss function considered guides the prediction method.

Some examples....

# SQUARED ERROR LOSS

Want to classify with  $f(x) = \mathbb{E}(Y|X = x)$  in order to minimize  $\text{EPE}(f)$ , but we can't calculate this classifier since we don't know  $p(x, y)$ .

## Least Squares

Least squares estimates  $f(x) = \mathbb{E}(Y|X = x)$  which minimizes  $\mathbb{E}_{(X,Y)}[L(Y, f(X))]$ , by minimizing observed or empirical  $L(Y, f(X))$  among all linear models.

## k-Nearest Neighbors

Nearest neighbors estimates  $f(x) = \mathbb{E}(Y|X = x)$  by using the mean of  $y_i$  values for  $x_i \in N_k(x)$  where  $N_k(x)$  is a small neighborhood of  $x$ .

# SQUARED ERROR LOSS

## Some Nice Properties of kNN

A larger sample means more observations close to  $x$  (i.e.  $N_k(x)$  is tighter) which produces a more stable  $\hat{f}(x)$  (given  $\mathbb{E}(Y|X)$  is continuous enough).

Actually, it is not difficult to show that if  $N, k \rightarrow \infty$  and  $k/N \rightarrow 0$ , then  $\hat{f}(x) \rightarrow \mathbb{E}(Y|X = x)$ .

Why not just use kNN, then? Often don't have large sample size, also curse of dimensionality which we address in a bit.

# ABSOLUTE LOSS

We replace the  $L_2$  metric by the  $L_1$  metric giving  $L(Y, f(X)) = |Y - f(X)|$ .

It can be shown that  $f(x) = \text{median}(Y|X = x)$  minimizes

$$\text{EPE}(f) = \mathbb{E}|Y - f(X)|.$$

As before, implies regression and kNN methods:

- ▶ Regression: least absolute deviations (LAD) regression or LAR (least absolute residuals).
- ▶ kNN: use median instead of average for calculating distance used to find each  $N_k(x)$ .

# CLASSIFICATION LOSS

## Definition

Want a function  $f(X)$  for predicting class  $Y$  given values of  $X$ . This theory requires a loss function:

$$L(Y, f(X)) : [K] \times [K] \rightarrow [0, \infty)$$

for penalizing prediction errors.

## Multiple Classes

Some example loss functions:

1. Loss function defined based on a penalty matrix,  $\mathbf{L} = [L(k, \ell)]_{K \times K}$  for  $K$  classes.
2. Zero-one loss corresponding to the number of misclassifications.  
Defined as  $L(k, \ell) = 1$  if  $k \neq \ell$  and 0 otherwise.



# CLASSIFICATION LOSS

EPE

Then,

$$\text{EPE}(f) = \mathbb{E}_{(X,Y)}[L(Y, f(X))] = \mathbb{E}_X \sum_{k=1}^K L(k, f(X)) \Pr(k|X).$$

With 0-1 loss, it can be shown that

$$\hat{f}(x) = k \quad \text{if} \quad \Pr(Y = k | X = x) = \max_{k' \in [K]} \Pr(Y = k' | X = x)$$

minimizes  $\text{EPE}(f)$ .

# BAYES CLASSIFIER

“Classify to the most probable class, using the conditional distribution  $\Pr(Y|X)$ ”

- ▶ Error rate of this classifier is called the Bayes rate.
- ▶ Bayes classifier is best under zero-one loss.
- ▶ In practice, can't calculate Bayes classifier because don't know  $p(x, y)$ .

kNN attempts to estimate the Bayes classifier by estimating  $\Pr(Y = k|X = x)$  by  $\hat{\Pr}(Y = k|N_k(x))$ .

# Methods Using Discriminant Analysis

# LINEAR DISCRIMINANT ANALYSIS (LDA)

- ▶ LDA assumes a Gaussian mixture with common covariance matrix,  $\Sigma$ .
- ▶ Want to estimate  $\Pr(Y = k | X = x)$ . By Bayes' Rule:

$$\Pr(Y = k | X = x) \propto \Pr(X = x | Y = k) \Pr(Y = k).$$

- ▶ When comparing two classes, it is sufficient to look at the log-ratio:

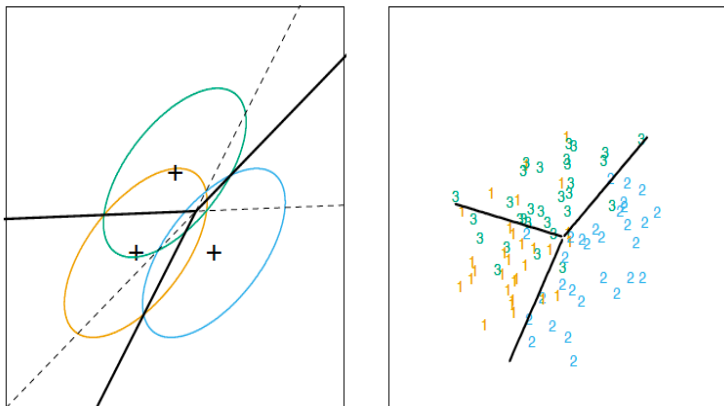
$$\begin{aligned} \log \frac{\Pr(Y = k | X = x)}{\Pr(Y = \ell | X = x)} &= \log \frac{\Pr(G = k)}{\Pr(G = \ell)} - \frac{1}{2}(\mu_k - \mu_\ell)^T \Sigma^{-1}(\mu_k - \mu_\ell) \\ &\quad + x^T \Sigma^{-1}(\mu_k - \mu_\ell) \end{aligned}$$

- ▶ Above implies, that the decision boundary between classes  $k$  and  $\ell$  is linear in  $x$ ; in  $p$  dimensions a hyperplane
- ▶ The linear discriminant function is then

$$\delta_k(x) = \log \Pr(Y = k) - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + x^T \Sigma^{-1} \mu_k$$

and  $\hat{f}(x) = \arg \max_k \delta_k(x)$ .

- ▶ Natural ways to estimate  $\hat{\pi}_k$  (with  $N_k/N$ ),  $\hat{\mu}_k$  (with  $\sum_{y_i=k} x_i/N_k$ ), and  $\hat{\Sigma}$  (with the pooled estimate of the covariance matrix).



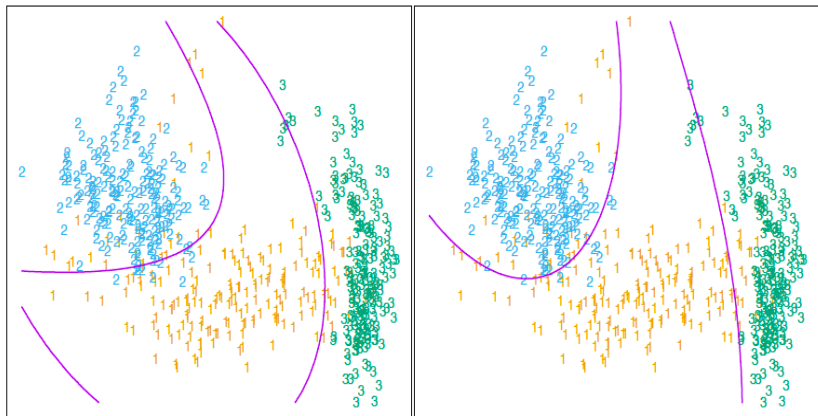
**FIGURE 4.5.** The left panel shows three Gaussian distributions, with the same covariance and different means. Included are the contours of constant density enclosing 95% of the probability in each case. The Bayes decision boundaries between each pair of classes are shown (broken straight lines), and the Bayes decision boundaries separating all three classes are the thicker solid lines (a subset of the former). On the right we see a sample of 30 drawn from each Gaussian distribution, and the fitted LDA decision boundaries.

# QUADRATIC DISCRIMINANT FUNCTIONS

If we assume different covariance matrices  $\Sigma_K$  for the classes  $k = 1, \dots, K$ , the discriminant function based on log-likelihood-ratio will be quadratic.

$$\delta_k(x) = \log Pr(Y = k) - \frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k)$$

Thus, the discriminant boundary ( $\{x : \delta_k(x) = \delta_\ell(x)\}$ ) is a quadratic function.

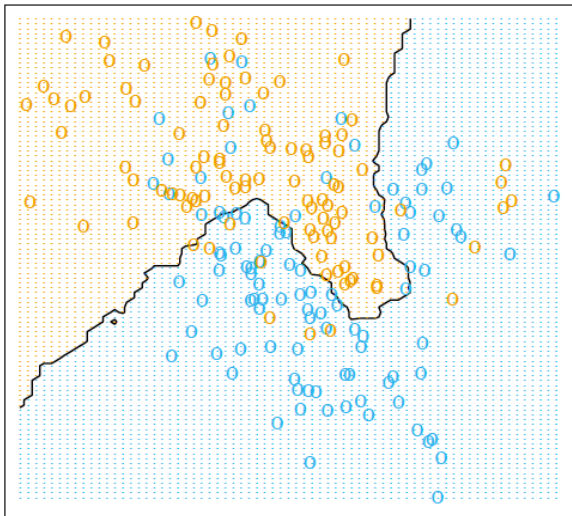


**FIGURE 4.6.** *Two methods for fitting quadratic boundaries. The left plot shows the quadratic decision boundaries for the data in Figure 4.1 (obtained using LDA in the five-dimensional space  $X_1, X_2, X_1X_2, X_1^2, X_2^2$ ). The right plot shows the quadratic decision boundaries found by QDA. The differences are small, as is usually the case.*

# METHOD OF NEAREST NEIGHBORS

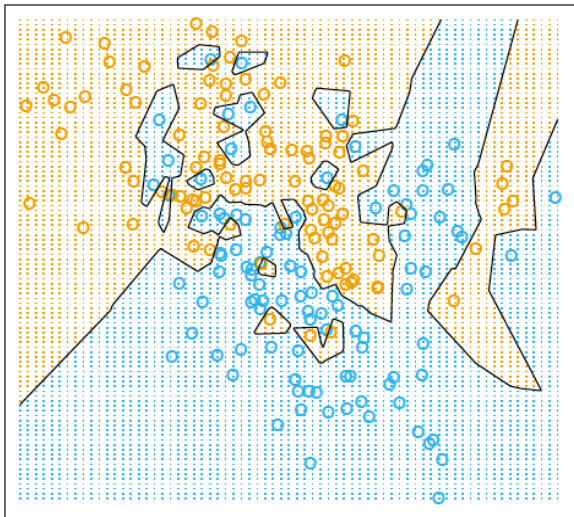
- ▶ Define  $N_k(x)$  as the neighborhood of  $x$  containing the  $k$  ‘nearest’ points in the training set.
- ▶ A distance metric is (implicitly) needed for kNN methods – most popular choice is Euclidean distance.
- ▶ Prediction is then  $\hat{f}(x) = \frac{1}{k} \sum_{x_i \in N_k(x)} y_i$ .
- ▶ Obviously residual sum of squares minimized at  $k = 1$  (all training classifications are correct).
- ▶ Models fit using kNN is less rigid than linear methods like LDA and generate non-linear prediction functions.



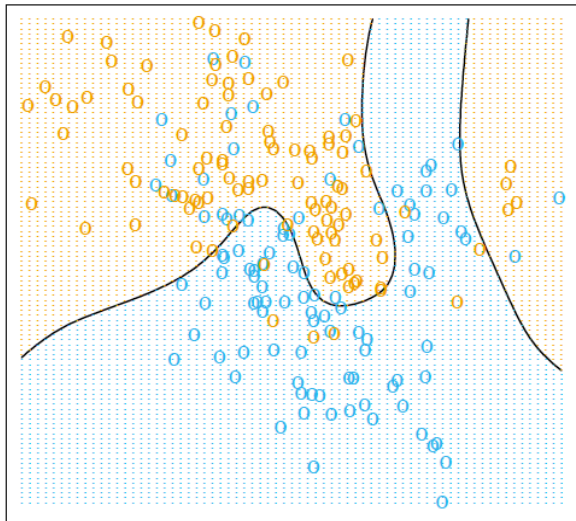


**FIGURE 2.2.** The same classification example in two dimensions as in Figure 2.1. The classes are coded as a binary variable (BLUE = 0, ORANGE = 1) and then fit by 15-nearest-neighbor averaging as in (2.8). The predicted class is hence chosen by majority vote amongst the 15-nearest neighbors.

### 1-Nearest Neighbor Classifier



**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.*



**FIGURE 2.5.** *The optimal Bayes decision boundary for the simulation example of Figures 2.1, 2.2 and 2.3. Since the generating density is known for each class, this boundary can be calculated exactly (Exercise 2.2).*

# PROPERTIES

## Linear methods

Low variance, high bias

- ▶ Unbiased if true model is linear, biased if true model is non-linear.
- ▶ Stable (i.e., individual observations not very influential) and smooth.
- ▶ Nice analytical results.

## Nearest Neighbor methods

High variance, low bias

- ▶ No assumption on the model.
- ▶ High variance in the predictions since each  $\hat{Y}$  is calculated only by a few observations.

# VARIANTS OF THE TWO METHODS

- ▶ Kernel methods can be applied to distance metrics. E.g. smoother kernel to replace the kNN “kernel”.
- ▶ Local weights, or weights varying across dimensions, to make linear methods less rigid.
- ▶ Linear models fit to a basis expansion of the original inputs: this expands the class of models considered.

# Curse of Dimensionality

# ASSUMPTIONS MADE IN CLASSIFICATION

## Assumptions

- ▶ Probability of class labels are continuous over feature values.
- ▶ The distance metric or kernel function is meaningful for the classification problem.
- ▶ Test sample will be drawn from the same distributions as the training sample,  $p(x, y)$ .

If you have infinite data and unlimited computational power and storage, classification is easy. Often not the case...

- ▶ For finite-size training sample, don't have enough observations to make predictions everywhere.
- ▶ Bayes rate can tell us about theoretical limits of performance.

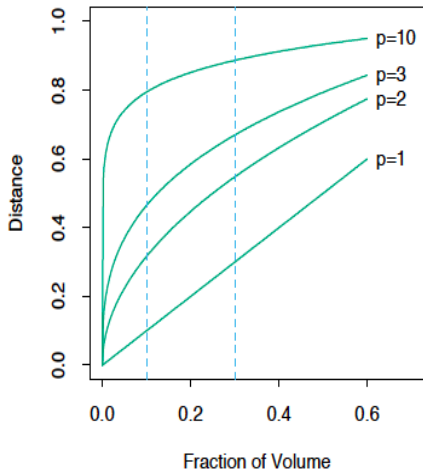
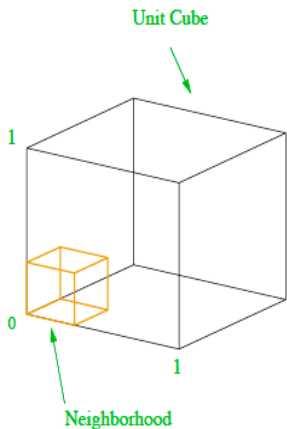
# HIGH DIMENSIONS

Recall:  $p$  is the dimension of the input space.

Assume inputs are uniformly distributed in a  $p$ -dimensional unit cube.

- ▶ Suppose we construct a hypercubical neighborhood about a target point to capture a fraction  $r$  of the observations. This corresponds to a fraction  $r$  of the unit volume, so the expected edge length is  $e_p(r) = r^{1/p}$ .
- ▶ In ten dimensions  $e_{10}(0.01) = 0.63$  and  $e_{10}(0.1) = 0.80$ , while the entire range is only 1.0. So to capture 1% or 10% of the data, we must cover 63% or 80% of the range of each input.
- ▶ Such neighborhoods are no longer “local.”





**FIGURE 2.6.** The curse of dimensionality is well illustrated by a subcubical neighborhood for uniform data in a unit cube. The figure on the right shows the side-length of the subcube needed to capture a fraction  $r$  of the volume of the data, for different dimensions  $p$ . In ten dimensions we need to cover 80% of the range of each coordinate to capture 10% of the data.

# CURSE OF DIMENSIONALITY

Recall:  $p$  is the dimension of the input space.

Assume inputs are uniformly distributed in a  $p$ -dimensional unit cube.

- ▶ For a random point  $\{x_1, \dots, x_p\}$ ,  $x_i \sim \text{Unif}(0, 1)$ , iid.
- ▶ It can be shown that  $\mathbb{E} \min(x_i) = \frac{1}{p+1}$ .
- ▶ This implies that for any point, it is very close to at least one boundary. Inference at the boundary is usually difficult.

# VARIANCE-BIAS DECOMPOSITION

## Model

$$Y = f(X) + \epsilon \quad \text{with } \mathbb{E}[\epsilon] = 0 \quad \text{and } \text{Var}[\epsilon] = \sigma^2$$

Let's consider the expected prediction error (EPE) for a model  $\hat{f}$  using  $L_2$  loss at a single test point  $(x_0, y_0)$ :

$$\begin{aligned} \text{EPE}(\hat{f}) &= \mathbb{E}(y_0 - \hat{f}(x_0))^2 \\ &= \mathbb{E}(f(x_0) + \epsilon - \mathbb{E}(\hat{f}(x_0)) + \mathbb{E}(\hat{f}(x_0)) - \hat{f}(x_0))^2 \\ &= \mathbb{E}(\epsilon^2) + (f(x_0) - \mathbb{E}(\hat{f}(x_0)))^2 + \text{var}(\hat{f}(x_0)) \\ &= \text{Irreducible error} + \text{Bias}^2 + \text{var}(\hat{f}(x_0)) \end{aligned}$$

where

$$\text{Irreducible error} = \sigma^2,$$

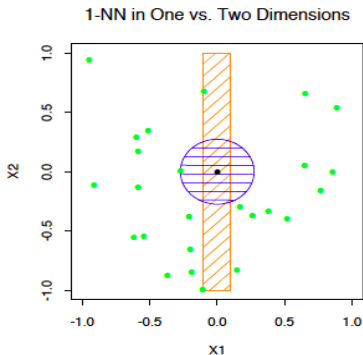
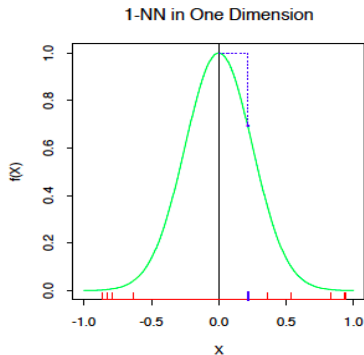
and

$$\text{Bias} = f(x_0) - \mathbb{E}(\hat{f}(x_0)).$$

# HOW DIMENSIONS AFFECT ESTIMATION

True relationship:  $Y = e^{-8\|X\|^2}$  (no measurement error).

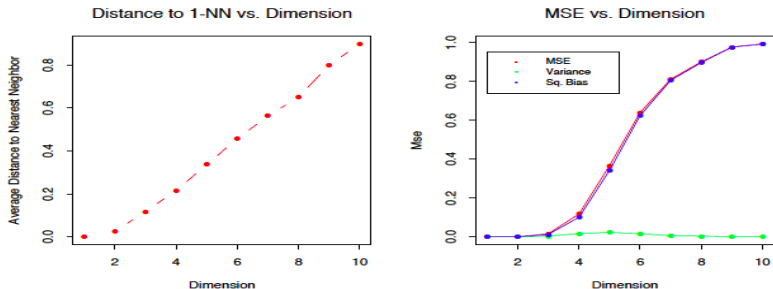
Use nearest-neighbor to estimate  $Y$  at  $X = 0$ .



**FIGURE 2.7.** A simulation example, demonstrating the curse of dimensionality and its effect on MSE, bias and variance. The input features are uniformly distributed in  $[-1, 1]^p$  for  $p = 1, \dots, 10$ . The top left panel shows the target function (no noise) in  $\mathbb{R}$ :  $f(X) = e^{-8\|X\|^2}$ , and demonstrates the error that 1-nearest neighbor makes in estimating  $f(0)$ . The training point is indicated by the blue tick mark. The top right panel illustrates why the radius of the 1-nearest neighborhood increases with dimension  $p$ . The lower left panel shows the average radius of the 1-nearest neighborhoods. The lower-right panel shows the MSE, squared bias and variance curves as a function of dimension  $p$ .

# HOW DIMENSIONS AFFECT ESTIMATION

- ▶ As the number of dimensions increase, the distance of the nearest neighbor to  $X = 0$  increases.
- ▶ The nearest neighbor estimate is therefore down-biased.
- ▶ The function  $Y = f(X)$  is symmetric about different dimensions of  $X$ .
- ▶ Actually, it only depends on the distance between the nearest neighbor at  $X = 0$ .
- ▶ The variability of the estimate is then decided by the variability of the distance to NN.



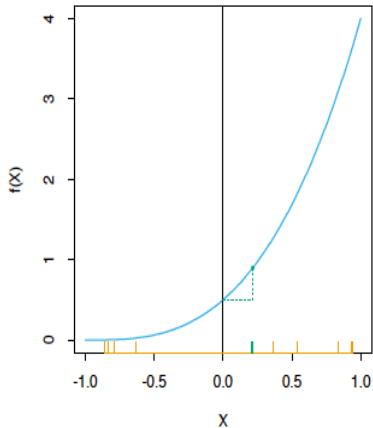
**FIGURE 2.7.** A simulation example, demonstrating the curse of dimensionality and its effect on MSE, bias and variance. The input features are uniformly distributed in  $[-1, 1]^p$  for  $p = 1, \dots, 10$ . The top left panel shows the target function (no noise) in  $\mathbb{R}$ :  $f(X) = e^{-8\|X\|^2}$ , and demonstrates the error that 1-nearest neighbor makes in estimating  $f(0)$ . The training point is indicated by the blue tick mark. The top right panel illustrates why the radius of the 1-nearest neighborhood increases with dimension  $p$ . The lower left panel shows the average radius of the 1-nearest neighborhoods. The lower-right panel shows the MSE, squared bias and variance curves as a function of dimension  $p$ .

# HOW DIMENSIONS AFFECT ESTIMATION

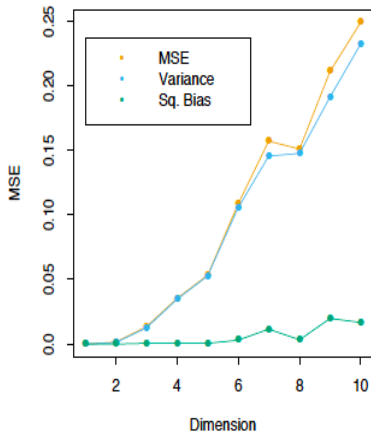
- ▶ Consider another case:  $Y = \frac{1}{2}(X_1 + 1)^3$ .
- ▶ The function depends on only one dimension, i.e., the other dimensions are irrelevant for learning this function.
- ▶ This function doesn't peak at 0 and therefore the bias isn't as prominent.
- ▶ Variability of the estimate depends on distance to NN along  $X_1$ , which increases as the number of irrelevant dimensions increases.



1-NN in One Dimension



MSE vs. Dimension



**FIGURE 2.8.** A simulation example with the same setup as in Figure 2.7. Here the function is constant in all but one dimension:  $F(X) = \frac{1}{2}(X_1 + 1)^3$ . The variance dominates.

# Statistical Models

# PREDICTIVE RELATION

## How We Build Predictors

Predictive relation between  $Y$  and  $X$  depends on the definition of “goodness of fit”, usually determined a loss function.

Examples:

- ▶ **L<sub>2</sub> Loss:** Best Estimate  $f(x) = \mathbb{E}(Y|X = x)$ .
- ▶ **L<sub>1</sub> Loss:** Best Estimate  $f(x) = \text{median}(Y|X = x)$ .

## kNN

- ▶ Nearest neighbor can be viewed as local direct estimates of  $f(x)$ .
- ▶ BUT, nearest neighbor methods run into trouble when the dimension of the input space becomes large.
- ▶ Moreover, if the relation between  $Y$  and  $X$  is known to be more structured, kNN methods aren't optimal.

# PREDICTIVE RELATION

For functions  $f$  and  $g$  that satisfy

$$f(x_i) = g(x_i), \quad i = 1, \dots, n,$$

their fit to the observed data  $(x_i, y_i); i = 1, \dots, n$  is the same.

The above fact leads to the definition of some kinds of equivalent models.

## Identifiability

- ▶ By constraining the model family, the hope is that within the model class, there are no equivalent models].
- ▶ When this is not the case, can have identifiability issues. Which model –  $f$  or  $g$  in the above definition – do we use?

## Occam's Razor

General belief: the more complicated a model, the more likely it is to give predictions far away from the truth at points not close to observed  $x_i$ 's.

- ▶ Complexity of models usually controlled with constraints.
- ▶ Generally require that the estimated model exhibit some kind of regular behavior in small neighborhoods of the input space.
- ▶ Modeling usually carried out using structured model families, basis expansions, and kernel/local regression.

How do we determine the model parameters?  
(multiplier of penalty term, width of kernel, number of basis functions, ...)