

Lecture 3

- K-NN vs Linear Regression
- Decision Theory for Classification
- The Bayes Classifier
- K-NN and Lin. Regression for Classific.
- The Bias-Variance Trade-Off

There is no free lunch in statistical learning: (12)
no one method dominates all other methods over all possible data sets. Usually, the more we assume, the more we can learn from data and the more accurate our predictions are, provided our assumptions are correct.

Last time we discussed two quite different regression methods:

- ① The k-NN method: $f(x) \approx \hat{f}(x) = \frac{1}{k} \sum_{i: x_i \in N_k(x)} y_i$
 - "Flexible" method: it approximates $f(x)$ by a locally constant function (class of these functions is large)
non-parametric
 - Has a nice theoretical property: $\hat{f}(x) \rightarrow f(x) = E[Y|X=x]$ as $N, k \rightarrow \infty$ and $k/N \rightarrow 0$ large
 - But it suffers from the curse of dimensionality. Prediction errors are large if $p \gg 1$. regression function, which is the L_2 -optimal prediction rule.
- ② The linear regression method: $f(x) \approx \hat{f}(x) = x^T \hat{\beta}$, $\hat{\beta} = \underbrace{(X^T X)^{-1} X^T y}_{\text{based on the training data.}}$
 - "Rigid" method: it approximates $f(x)$ by a globally linear function (class of these functions is small)
parametric
 - If the true relationship between the output Y and input X is approximately linear, $Y \approx f(X) = X^T \beta$, then the method can overcome the curse of dimensionality.
 - But if the linearity assumption is wrong \Rightarrow predictions are very bad.

Many popular regression methods are enhanced modifications of these two methods. There is a whole spectrum of methods between the rigid linear regression and the extremely flexible 1-nearest-neighbor method.

Here are some immediate ideas on how to enhance the k-NN and linear methods:

- The k-NN approximation can be written as follows: $\hat{f}(x) = \frac{1}{k} \sum_{i=1}^k w_i y_i$, where
Instead of 0/1 weights, we can have weights w_i that decrease smoothly to zero with distance from the target input x : $\hat{f}(x) = \frac{1}{N} \sum_{i=1}^N \underbrace{w_i(\|x - x_i\|)}_{\text{kernel}} y_i$
This leads to the so called kernel methods. must be chosen wisely in high dim.
 $w_i = \begin{cases} 1, & \text{if } x_i \in N_k(x) \\ 0, & \text{if } x_i \notin N_k(x) \end{cases}$
- We can extend the standard linear regression by allowing non-linear functions of each input, while maintaining additivity: $f(x) = \sum_{j=1}^p f_j(x_j)$.
This leads to the generalized additive models. arbitrary function

The k -NN and linear regression methods can be also used for solving classification problems. They naturally appear under the framework of Statistical Decision Theory for classification (for regression, see Lecture 1, p.4-5)

Let $X \in \mathbb{R}^p$ be a random input and $G \in \mathcal{G} = \{g_1, \dots, g_K\}$ be the corresponding random qualitative output. We want to find a function (prediction rule) $f: \mathbb{R}^p \rightarrow \mathcal{G}$ such that $\hat{G} = f(X) \approx G$ is an accurate prediction.

Remark: Re notation: It is a bit heavy, but $\mathcal{G} = \{1, \dots, K\}$ has a drawback: it assumes linear order between classes, which in general does not exist, e.g., $\mathcal{G} = \{\text{"stroke"}, \text{"drug overdose"}, \text{"epileptic seizure"}\}$

To measure the accuracy of prediction (goodness of f), we need a loss function $L: \mathcal{G} \times \mathcal{G} \rightarrow \mathbb{R}$.

If the expected loss is small, f is good; otherwise f is bad:

$$E[L(G, \hat{G})] = E[L(G, f(X))] = \begin{cases} \text{small} \Rightarrow f \text{ is good,} \\ \text{large} \Rightarrow f \text{ is bad.} \end{cases}$$

The most popular and often used loss function is the zero-one loss, where all misclassifications are charged a single unit:

$$L(G, \hat{G}) = I(G \neq \hat{G}) = \begin{cases} 0 & \text{if } G = \hat{G} \\ 1 & \text{if } G \neq \hat{G} \end{cases}$$

indicator variable

Remark: In regression settings, $L: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ and we used the squared error loss $L(Y, \hat{Y}) = (Y - \hat{Y})^2$. In classification settings, it does not make sense.

The corresponding expected loss, the expected prediction error (EPE), is then:

$$EPE(f) = E[I(G \neq \hat{G})] = E[I(G \neq f(X))] \leftarrow \text{this expected value is w.r.t. the joint distribution of } X \text{ and } G.$$

Our goal is to find $f(x)$ that minimizes $EPE(f)$.

As in the regression case, we condition on X and use the law of total expectation:

$$EPE(f) = E[E[I(G \neq f(X)) | X]] \rightarrow \min$$

wrt X wrt G

to minimize $EPE(f)$, we need to find $f(x)$ that minimizes this expression

The inner expectation: $E[I(G \neq f(x)) | X=x] = \sum_{k=1}^K I(g_k \neq f(x)) \cdot P(G=g_k | X=x)$

$$\Rightarrow f(x) = \arg \min_{g_s \in \mathcal{G}} \sum_{k=1}^K I(g_k \neq g_s) P(G=g_k | X=x) = \arg \min_{g_s \in \mathcal{G}} \sum_{k \neq s} P(G=g_k | X=x)$$

$= \begin{cases} 0, & k=s \\ 1, & k \neq s \end{cases}$

So, the prediction rule that minimizes $EPE(f)$ is

$$f(x) = \arg \max_{g \in \mathcal{G}} P(G=g | X=x)$$

This solution is called the Bayes classifier.

$$= \arg \min_{g_s \in \mathcal{G}} (1 - P(G=g_s | X=x))$$

since $\sum_{k=1}^K P(G=g_k | X=x) = 1$

$$= \arg \max_{g_s \in \mathcal{G}} P(G=g_s | X=x)$$

If $X=x$, then the best prediction is $\hat{G} = g_k \Leftrightarrow P(G=g_k | X=x) = \max_{g \in \mathcal{G}} P(G=g | X=x)$

zero-one sense

Remark: The Bayes classifier $f(x) = \arg \max_{g \in G} P(G=g | X=x)$ is

a classification analog of the regression function $f(x) = E[Y | X=x]$

- Both are obtained by minimizing the EPE, but w.r.t. different loss functions.
- Both functions are unknown and we want to learn them from the data.
- The Bayes classifier and the regression function are, in fact, related.

Consider a 2-class problem: $G \in G = \{g_0, g_1\}$.

Define a binary dummy-variable: $Y=0$ if $G=g_0$ and $Y=1$ if $G=g_1$.

$$\Rightarrow E[Y | X=x] = P(G=g_1 | X=x)$$

$$\Rightarrow \arg \max_{g \in G} P(G=g | X=x) = \begin{cases} g_0 & \text{if } E[Y | X=x] < 1/2 \\ g_1 & \text{if } E[Y | X=x] > 1/2 \end{cases}$$

Remark This can be generalized to a K-class problem via

$$Y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_K \end{bmatrix} \quad Y_i = \begin{cases} 1 & \text{if } G=g_i \\ 0 & \text{if } G \neq g_i \end{cases}$$

In applications, the conditional probability $P(G=g | X=x)$ is unknown, and, therefore, we can't use the Bayes classifier directly.

One strategy is to estimate $P(G=g | X=x)$ from the data, and then classify a given observation $X=x$ to the class with highest estimated probability. One such method is

- The k-NN Classifier ← very much in the spirit of its regression brother.

1. Choose $k \in \{1, \dots, N\}$, where N is the sample size of training data $\{(x_1, g_1), \dots, (x_N, g_N)\}$

2. Relax the condition $X=x$ by $X \in N_k(x)$, where $N_k(x)$ is the set of the k closest to x inputs x_i in the data.

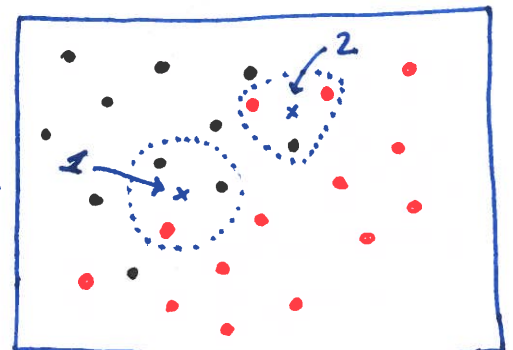
3. Estimate the conditional probability $P(G=g | X=x)$ by the fraction of points in $N_k(x)$ whose class is g : $P(G=g | X=x) \approx \frac{1}{k} \sum_{x_i \in N_k(x)} I(g_i=g)$

4. Approximate the Bayes classifier:

$$f(x) \approx \hat{f}(x) = \arg \max_{g \in G} \underbrace{\frac{1}{k} \sum_{x_i \in N_k(x)} I(g_i=g)}_{\text{not important}} = \arg \max_{g \in G} \sum_{x_i \in N_k(x)} I(g_i=g)$$

So, the k-NN classifier is simply a majority vote in the neighborhood $N_k(x)$.

Despite the fact that the k-NN classifier is a very simple approach, it often performs very well (close to the optimal Bayes classifier) x_2 in low dimensions. In high dimensions, however, it suffers from the curse of dimensionality.



Remark: The choice of k has a drastic effect on the k-NN classifier. We will discuss later how to choose k in a systematic way.

• class 1
• class 2
 $k=3$
 $p=2$

Another strategy is to model the Bayes classifier instead of estimating it. One method that implements this idea is the linear regression model.

• Linear Regression for Classification

Consider a 2-class problem: $G \in \mathcal{G} = \{g_0, g_1\}$.

Define a binary dummy-variable: $Y=0$ if $G=g_0$ and $Y=1$ if $G=g_1$.

Convert the training data: $\{(x_1, g_1), \dots, (x_N, g_N)\} \rightsquigarrow \{(x_1, y_1), \dots, (x_N, y_N)\}$, where

Use linear regression for predicting output Y from input X :

$$\hat{Y} = \hat{f}(X) = X^T \hat{\beta}, \text{ where } \hat{\beta} = (X^T X)^{-1} X^T y,$$

$$y_i = \begin{cases} 0 & \text{if } g_i = g_0 \\ 1 & \text{if } g_i = g_1 \end{cases}$$

$$X = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_p \end{bmatrix} \text{ and } \mathbf{X} = \begin{bmatrix} 1 & x_{11} & \dots & x_{1p} \\ \vdots & \vdots & & \vdots \\ 1 & x_{N1} & \dots & x_{Np} \end{bmatrix} = \begin{bmatrix} x_1^T \\ \vdots \\ x_N^T \end{bmatrix} \text{ and } y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

Finally, predict class G of input X as follows:

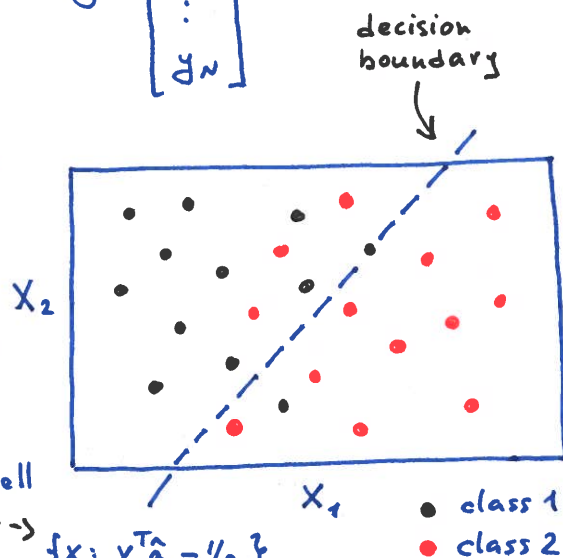
$$\hat{G} = \begin{cases} g_0 & \text{if } \hat{Y} < 1/2 \\ g_1 & \text{if } \hat{Y} > 1/2 \end{cases}$$

The sets of inputs in \mathbb{R}^p classified as g_0 and g_1 correspond to

$$\{X: X^T \hat{\beta} < \frac{1}{2}\} \leftarrow g_0$$

$$\{X: X^T \hat{\beta} > \frac{1}{2}\} \leftarrow g_1$$

This approach will work well if \exists a hyperplane $\dots \rightarrow \{X: X^T \hat{\beta} = 1/2\}$ that approximately divides the two classes.



Remark: This can be generalized to the K -class problem using dummy-variable

$$Y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_K \end{bmatrix} \text{ where } Y_i = \begin{cases} 0 & \text{if } G \neq g_i \\ 1 & \text{if } G = g_i \end{cases}$$

Use multiple linear regression
 \downarrow will discuss later

for predicting Y from X : $\hat{Y}^T = X^T \hat{B} \Rightarrow$ predict class G by $\hat{G} = g_{\arg \max_k \hat{Y}_k}$
 We will discuss this in more detail when we will be discussing classification problems.

As mentioned at the beginning, no one method dominates all others over all possible data sets. On a particular data set, however, one specific method may perform better than others. Selecting the best method from a collection of methods for a given training data set is called model selection, and it is one of the most difficult tasks in statistical learning.

Example: • Should we use a more flexible k -NN method or a less flexible linear regression?
 • If k -NN, then what value of k should we use?
 • If linear regression, then how many parameters $\beta_0 \dots \beta_p$ should we use?
 (May be some inputs X_i are not relevant for predicting Y and we can set $\beta_i = 0$)

We will discuss model selection later. Here let's discuss two competing statistical forces that make model selection a challenging task.

The Bias-Variance Trade-Off

16

Let's consider the regression setting with the squared error loss function.

Remark: For other reasonable loss functions and for classification, the story is similar.

Suppose that: $Y = f(X) + \varepsilon$, where $IE[\varepsilon] = 0$, $W[\varepsilon] = \sigma^2$

Then $IE[Y|X=x] = IE[f(x) + \varepsilon|X=x] = f(x)$ ← the regression function that we want to learn. not necessarily small.

Let $\mathcal{I} = \{(x_1, y_1), \dots, (x_N, y_N)\}$ be the training data,

and $\hat{f}(X) \approx f(X)$ be an approximation of $f(X)$ obtained from \mathcal{I} by some method

So, our prediction for output Y from input X is $\hat{Y} = \hat{f}(X)$.

A reasonable model \hat{f} should perform well on the training data \mathcal{I} .

This performance is quantified by the training error:

$$\bar{err} = \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{f}(x_i)) = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{f}(x_i))^2$$

Q: Can we use \bar{err} for choosing between different models \hat{f} ?

A: No. We are not really interested how well \hat{f} works on \mathcal{I} .

Instead, we are interested in the accuracy $Y \approx \hat{f}(X)$, where X is new previously

Example: Suppose, we want to use the k -NN method and need to select k . unseen input.

To minimize \bar{err} , regardless of the training data \mathcal{I} , we should always use $k=1$

In this case, $\hat{f}(x_i) = y_i$ (since $N_1(x_i) = x_i$), and $\bar{err} = 0$. which does not feel right especially in the context of the limiting result on p. 6:

So, instead of \bar{err} , we are interested in how well \hat{f} can "generalize" from the training data to new unseen data.

Let X be a new fixed input. The quantity of interest is the expected prediction error (EPE) at X , aka generalization or test error

$$Err(X) = IE[L(Y, \hat{f}(X))] = IE[(Y - \hat{f}(X))^2]$$

w.r.t. randomness in Y and \mathcal{I}

$$\hat{f} \rightarrow IE[Y|X=x]$$
$$k, N \rightarrow \infty \quad \frac{k}{N} \rightarrow 0$$

This test error can always be decomposed into the sum of three fundamental quantities:

$$Err(X) = IE[(\varepsilon + f(x) - \hat{f}(x))^2] = IE[\varepsilon^2] + IE[(f(x) - \hat{f}(x))^2] + 2 IE[\varepsilon \cdot (f(x) - \hat{f}(x))]$$

$$= \sigma^2 + IE[(f(x) - IE[\hat{f}(x)] + IE[\hat{f}(x)] - \hat{f}(x))^2]$$

$$= \sigma^2 + IE[(f(x) - IE[\hat{f}(x)])^2] + IE[(\hat{f}(x) - IE[\hat{f}(x)])^2]$$

$$+ 2 IE[(f(x) - IE[\hat{f}(x)])(IE[\hat{f}(x)] - \hat{f}(x))]$$

$$IE[(IE[\hat{f}(x)] - \hat{f}(x))] = 0$$

$$\underbrace{= 2 IE[\varepsilon] \cdot IE[f(x) - \hat{f}(x)]}_{\text{independent}} = 0$$

So
$$\text{Err}(X) = \underbrace{G^2}_{\text{irreducible error}} + \underbrace{(f(x) - \mathbb{E}[\hat{f}(x)])^2}_{\text{squared bias of } \hat{f}(x)} + \underbrace{\mathbb{V}[\hat{f}(x)]}_{\text{variance of } \hat{f}(x)}$$

← bias-variance trade-off.

All three terms are non-negative and we want them to be as small as possible.

- G^2 : irreducible error, we can't control it, even if we know the true value $f(x)$.
- $\mathbb{B}[\hat{f}(x)]^2$: squared bias of $\hat{f}(x)$, quantifies the error made by approximating the true value $f(x)$ by a model $\hat{f}(x)$. (IE is w.r.t to training data \mathcal{I})
- $\mathbb{V}[\hat{f}(x)]$: variance of $\hat{f}(x)$, quantifies the variability of $\hat{f}(x)$ w.r.t. \mathcal{I} . ideally, $\hat{f}(x)$ should not change a lot between different training sets.

can control to some degree

Generally: the more flexible the model \hat{f} is, the lower the bias $\mathbb{B}[\hat{f}(x)]^2$ but the higher the variance $\mathbb{V}[\hat{f}(x)]$

that is why "trade-off"

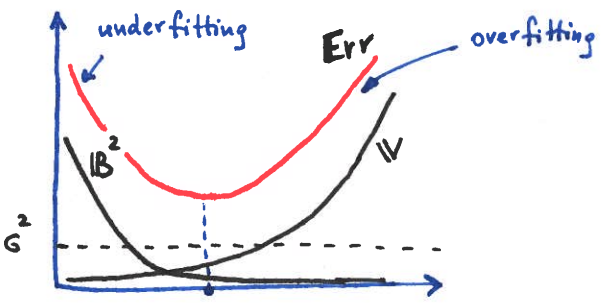
Remark: Intuitively, model flexibility

(aka model complexity) refers to the ability to fit many different functional forms. For instance, quadratic polynomials is a more flexible model than linear functions. Often (but not always), more flexible models have more parameters.

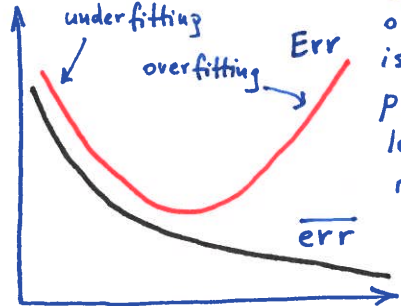
The concept of flexibility/complexity is formalized by the notion of the degrees

of freedom

Schematically, the picture looks as follows:



We want to choose the model flexibility to trade bias off with variance to minimize the test error.



The training error \bar{err} is a bad estimate of the test error: it does not account for model flexibility. The training error \bar{err} consistently \downarrow as flexibility \uparrow . Typically, if model is flexible enough, then $\bar{err} = 0$.

Remark: A U-shape

of the test error Err is a fundamental property of statistical learning that holds regardless of the statistical method being used.

aka effective number of parameters.

Remark: In applications, f is unknown, and it is not possible to compute Err , \mathbb{B} , \mathbb{V} . But it is useful to keep the general picture in mind

There are 2 bad cases:

- Overfitting: \hat{f} has small \bar{err} , but large Err (small bias, but large variance). The model adapts itself too closely to \mathcal{I} (finds patterns specific to \mathcal{I} , but not to f), and will not generalize well (have large test error)
- Underfitting: \hat{f} has large Err , because variance is small, but bias is large. The model is not flexible enough to approximate f well.

Remark: Regardless of whether Overfitting/Underfitting occurs, we almost always expect $\bar{err} < \text{Err}$, since most learning methods seek to minimize \bar{err} (directly or indirectly)

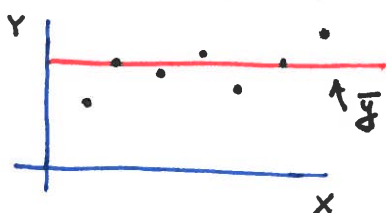
Example: To illustrate the bias-variance trade-off, let's consider the k -NN regression method:

$$f(x) \approx \hat{f}_k(x) = \frac{1}{k} \sum_{i: x_i \in N_k(x)} y_i, \quad k \in \{1, 2, \dots, N\}$$

- If $\underline{k=N} \Rightarrow \hat{f}_N(x) = \frac{1}{N} \sum_{i=1}^N y_i \leftarrow$ for a given \mathcal{T} , it is a constant.

In this case, we approximate $f(x)$ by a constant function.

effectively
one
parameter

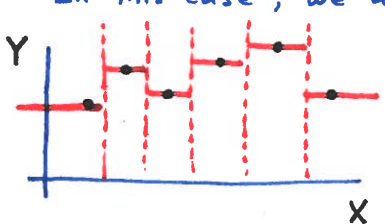


Whatever input x is,
our prediction for the output Y is $\hat{Y} = \bar{y}$
Very rigid (not flexible) method.

- If $\underline{k=1} \Rightarrow \hat{f}_1(x) = y_{i(x)}$, where $i(x) = \arg \min_{i=1, \dots, N} \text{dist}(x, x_i)$

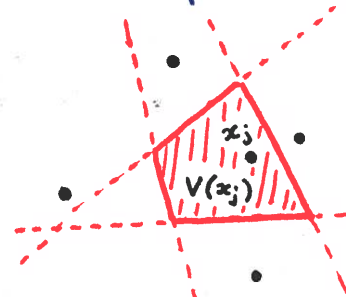
In this case, we approximate $f(x)$ by a locally constant function on a Voronoi tessellation of the \mathcal{T} .

effectively
 N
parameters



This is the most
flexible k -NN method.

If $p=2$



$$V(x_j) = \{x : x_j \text{ is the closest to } x \text{ among all } x_1, \dots, x_N\}$$

$$\forall x \in V(x_j) \quad \hat{f}_1(x) = y_j$$

Remark: As discussed above,
 $\text{err}(\hat{f}_1) = 0$

In general, k controls the flexibility of the k -NN model: $k \uparrow \Rightarrow \text{flexibility} \downarrow$

Let's compute the bias and the variance and check how they behave w.r.t. k .

We assumed that $Y = f(X) + \varepsilon \Rightarrow y_i = f(x_i) + \varepsilon_i$ ($\varepsilon_1, \dots, \varepsilon_N$ are iid $E[\varepsilon_i] = 0$
 $V[\varepsilon_i] = \sigma^2$)

Therefore:
$$\hat{f}_k(x) = \frac{1}{k} \sum_{i: x_i \in N_k(x)} (f(x_i) + \varepsilon_i)$$

Assume for simplicity that x_1, \dots, x_N in \mathcal{T} are fixed (nonrandom) and randomness arises from

$$B[\hat{f}_k(x)] = (f(x) - E[\hat{f}_k(x)])^2 = \left(f(x) - \frac{1}{k} \sum_{x_i \in N_k(x)} f(x_i)\right)^2 \quad y_1, \dots, y_N \text{ (from } \varepsilon_1, \dots, \varepsilon_N)$$

If flexibility $\uparrow \Rightarrow k \downarrow \Rightarrow B[\hat{f}_k(x)] \downarrow$ if $f(x)$ is reasonably smooth

(the closest neighbor is the most relevant for approximating $f(x)$, 2nd closest is a bit less relevant etc)

$$V[\hat{f}_k(x)] = V\left[\underbrace{\frac{1}{k} \sum f(x_i)}_{\text{constant}} + \frac{1}{k} \sum \varepsilon_i\right] = \frac{1}{k^2} \sum V[\varepsilon_i] = \frac{\sigma^2}{k}$$

If flexibility $\uparrow \Rightarrow k \downarrow \Rightarrow V[\hat{f}_k(x)] \uparrow$

Final Remark: For both regression and classification, choosing the correct level of flexibility is critical for the success of any learning method.