

Lecture 2

- Nearest-Neighbor Methods
- Linear Regression Model
- Curse of Dimensionality

Last time we have shown that in the regression settings ($X \in \mathbb{R}^p$, $Y \in \mathbb{R}$), the optimal (in the MSE sense) prediction of the output to input $X=x$ is given by the regression function: $f(x) = \mathbb{E}[Y|X=x]$.

This function is unknown and our goal is to estimate it.

Nearest-Neighbor Methods provide the most straightforward way to estimate $f(x)$.

The main idea is to estimate $\mathbb{E}[Y|X=x]$ "nonparametrically" directly from the data.

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

Then we can estimate the expectation by the empirical average: $\mathbb{E}[Y|X=x] \approx \text{Ave}(y_i | x_i = x)$

The problem with this naive approach is that typically there is at most one observation at any $x = x_i$, that is the set $\{y_i | x_i = x\}$ consists of 1 point or empty.

This difficulty can be resolved by relaxing condition $x_i = x$

Let's replace $x_i = x$ with $x_i \in N_k(x)$, where $N_k(x)$ is the set of the k closest inputs x_i in the data

This leads to the following approximation of the regression function:

$$(\#) \quad f(x) \approx \hat{f}(x) = \frac{1}{k} \sum_{i: x_i \in N_k(x)} y_i$$

k-nearest neighbor (k-NN method)

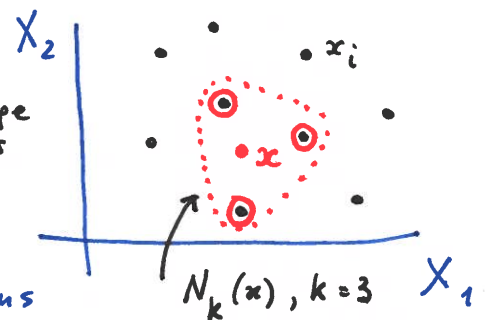
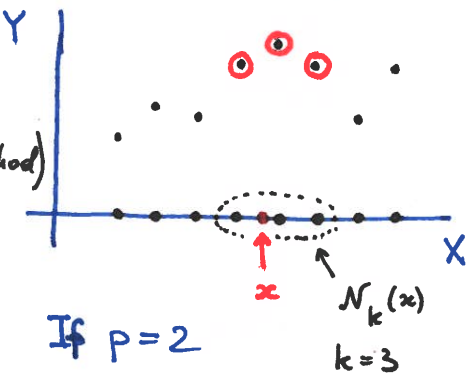
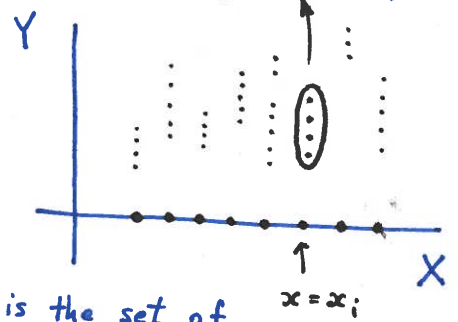
Two approximations are happening here:

1. Expectation is approximated by the empirical average.

If $N, k \rightarrow \infty$, then approximation becomes better thanks to the LLN

2. Condition $x_i = x$ is relaxed by $x_i \in N_k(x)$.

As $N \rightarrow \infty$, points in $N_k(x)$ become closer to x , and the corresponding approximation becomes better.



Important Fact: Under certain mild regularity conditions on the joint distribution of X and Y , it can be shown that

Devroy et al (1996)
"A probabilistic theory of pattern recognition."

$$(\#) \quad \hat{f}(x) \rightarrow f(x) = \mathbb{E}[Y|X=x] \quad \text{as } N, k \rightarrow \infty \text{ such that } \frac{k}{N} \rightarrow 0.$$

From a theoretical point of view, the k -NN method looks very good and it seems we can use it as a universal solution to the regression problem. But from a practical point of view, it has one drawback: the convergence (*) is a limiting result that holds when $N, k \rightarrow \infty$ (and $k/N \rightarrow 0$), but for finite N and k the approximation $f(x) \approx \hat{f}(x)$ in (#) may not be (and often is not) accurate. This is especially the case if $p = \dim X$ is large: the k -NN method suffers from the curse of dimensionality.

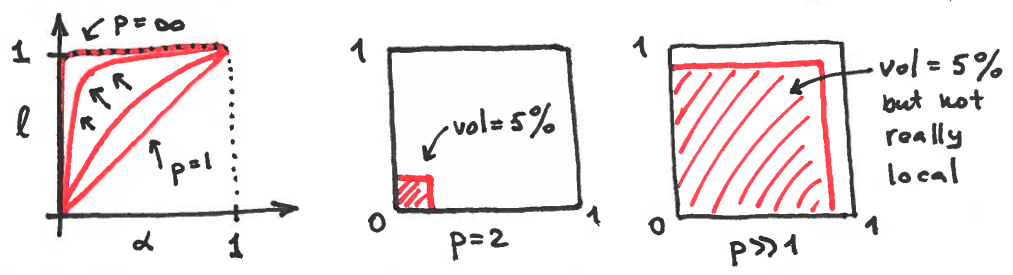
refers to phenomena that arise when analyzing data in high-dimensional spaces that do not occur in low-dim. spaces.

the convergence (*) still holds, but the rate of converges \downarrow as $p \uparrow$.

Let's consider several manifestations of the curse of dimensionality.

① The k -NN method is a local method: $f(x)$ is approximated based on the information in $N_k(x)$ that contains a fraction $\alpha = k/N$ of training sample. But in high dimensions the neighborhood $N_k(x)$ is not really "local". To illustrate this, consider inputs x_1, \dots, x_N uniformly distributed in $C_p = [0, 1]^p$. To capture a fraction α of the sample by a hypercubical neighborhood $C_p(l) = [0, l]^p$, we need to choose l such that $\text{vol}(C_p(l)) = \alpha$.

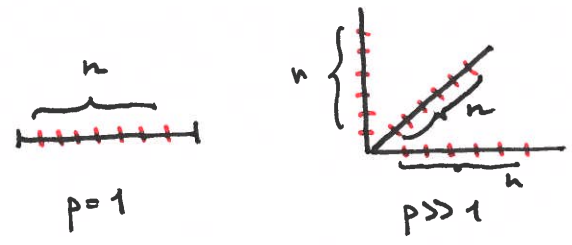
$\Rightarrow l^p = \alpha \iff l = \alpha^{1/p} \quad (\rightarrow 1 \text{ as } p \rightarrow \infty)$



Remark: For example, if $p = 20$ and $\alpha = 0.01$, then $l \approx 0.8$. That is to capture 1% of the data to form a local average, we must cover about 80% of the range of each input variable in 20-dim space.

②. Realistic training samples x_1, \dots, x_N sparsely populate the high-dimensional input spaces.

For example, if n measurements of $X_1 \in [0, 1]$ represent a dense sample for a single input problem, then to get the same "resolution" (sampling density) of C_p , we need to use the sample size $N = \underbrace{n \times \dots \times n}_p = n^p$, which is not feasible.



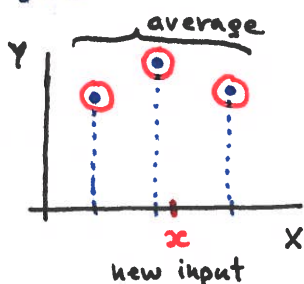
The problem is that the volume increases exponentially with adding extra dimensions.

for instance if $n=100$ and $p=10$
 \Downarrow
 $N = 100^{10}$

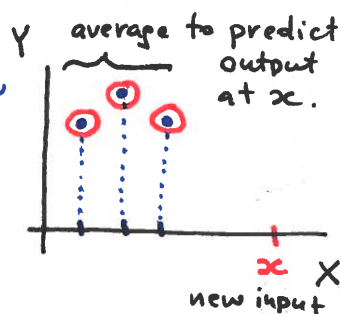


- ③ An important consequence of the sparse sampling in high dimensions is that all sample points are close to the "boundary" of the sample. This is a problem for the k -NN method: instead of interpolating between nearest neighbors $x_i \in N_k(x)$ (as (#) suggests), we will have to extrapolate from them.

Schematically: instead of nice picture suggested by our low-dim. intuition



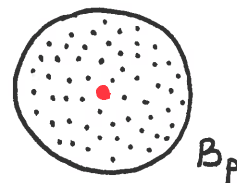
in high-dim, we will have the following picture



To illustrate this, consider N inputs uniformly distributed in B_p ← p -dimensional unit ball centered at the origin.

$$X_1, \dots, X_N \sim \mathcal{U}(B_p)$$

Low-dimensional intuition suggests that we have: Suppose we want to use a nearest-neighbor method to estimate the output at the origin.



Let D be the distance from the origin to the nearest neighbor (the closest X_i)

$$\Rightarrow D = \min \{ \|X_1\|, \dots, \|X_N\| \}$$

Norms $\|X_1\|, \dots, \|X_N\|$ are i.i.d. Let's find their CDF $F_{\|X\|}(d)$, $X \sim \mathcal{U}(B_p)$

$$F_{\|X\|}(d) = \mathbb{P}(\|X\| \leq d) = \frac{\text{vol}(B_p(d))}{\text{vol}(B_p)} = \frac{C \cdot d^p}{C \cdot 1^p} = d^p, \quad d \in [0, 1].$$

Let's find the CDF of D :

$$F_D(d) = \mathbb{P}(D \leq d) = 1 - \mathbb{P}(D > d)$$

$$= 1 - \mathbb{P}(\|X_1\| > d, \dots, \|X_N\| > d) = 1 - \prod_{i=1}^N \mathbb{P}(\|X_i\| > d)$$

$$= 1 - \prod_{i=1}^N (1 - \underbrace{\mathbb{P}(\|X_i\| \leq d)}_{F_{\|X\|}(d)}) = 1 - (1 - d^p)^N$$

p -dimensional ball of radius d centered at the origin.

Remark: C is a constant that depends on p :

$$C = \frac{\pi^{p/2}}{\Gamma(p/2 + 1)}$$

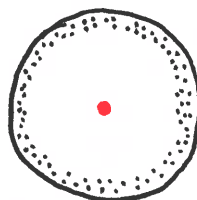
↑ Euler's gamma function

Let's use the median \tilde{d} as a proxy for a typical value of D .

$$F_D(\tilde{d}) = \frac{1}{2} \Rightarrow \tilde{d} = \left(1 - \left(\frac{1}{2}\right)^{1/N}\right)^{1/p} \rightarrow 1 \text{ as } p \rightarrow \infty.$$

\Rightarrow in high dimensions the picture looks like this:

For example, if $p = 20$, $N = 10^3 \Rightarrow \tilde{d} \approx 0.7$
 $p = 50$, $N = 10^3 \Rightarrow \tilde{d} \approx 0.86$



Remark: It is possible to derive an expression for $\tilde{d} = \mathbb{E}[D]$, but it is more complicated.

\Rightarrow using k -NN at the center leads to extrapolation

$B_p, p \gg 1$

So, if the sample size N is not extremely large and/or the dimension p of the input space is high, then the nearest neighbors $N_k(x)$ may not be close to the target input x , and this can result in large prediction errors.

This motivates us to look for alternative methods. Let's consider a completely different strategy: instead of estimating the regression function $f(x)$ directly from the data, let's model it, that is impose some structural assumptions on $f(x)$. ↘₂

Suppose $f(x)$ is linear: $f(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$

Let $\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}$ and $X = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_p \end{bmatrix} \Rightarrow f(x) = X^T \beta$ $X \in \mathbb{R}^{p+1}$
 $\beta \in \mathbb{R}^{p+1}$

Remark The linear model has been a mainstay of statistics for the last 50 years and remains one of the most important tools.

The mean squared error (the expected prediction error under the squared error loss function) is then:

Remark The term β_0 is called bias in ML.

$$MSE(f) = E[(Y - \hat{Y})^2] = E[(Y - X^T \beta)^2] = MSE(\beta) \leftarrow \text{function of } \beta.$$

To find the regression coefficients β that minimize MSE, we need to solve $\nabla_{\beta} MSE = 0$

$$\nabla_{\beta} MSE(\beta) = \nabla_{\beta} E[(Y - X^T \beta)^2] = E[\nabla_{\beta} (Y - X^T \beta)^2] = E[2(Y - X^T \beta) \nabla_{\beta} (Y - X^T \beta)] \quad \text{for } \beta.$$

↑ integration w.r.t. X and Y commutes with differentiation w.r.t. β .

$$= -2 E[(Y - X^T \beta) \nabla_{\beta} (X^T \beta)] = -2 E[(Y - X^T \beta) X] = 0.$$

$= X$, since the gradient is a vector (and vectors are columns)

This leads to: $E[YX] - E[\underbrace{X^T \beta X}_{\text{scalar}}] = 0 \Rightarrow E[XX^T \beta] = E[YX].$

Finally: $E[XX^T] \beta = E[YX] \Rightarrow \boxed{\beta = E[XX^T]^{-1} \cdot E[XY]} \quad (!)$ ↘₁

Remark:

$$XY = YX$$

since $Y \in \mathbb{R}$

Equation (!) gives theoretically optimal value of the regression parameters.

Let's now estimate β using the training data $(x_1, y_1), \dots, (x_N, y_N)$.

Let's estimate the expected values in (!) by the corresponding empirical averages.

$$\bullet E[XX^T] \approx \frac{1}{N} \sum_{i=1}^N x_i x_i^T = \frac{1}{N} \left(\overbrace{x_1 x_1^T}^{(p+1) \times (p+1)} + \dots + \overbrace{x_N x_N^T}^{(p+1) \times (p+1)} \right)$$

$$= \frac{1}{N} \begin{bmatrix} | & & | \\ x_1 & \dots & x_N \\ | & & | \end{bmatrix} \cdot \left(\begin{bmatrix} -x_1^T \\ 0 \end{bmatrix} + \dots + \begin{bmatrix} 0 \\ -x_N^T \end{bmatrix} \right)$$

$$= \frac{1}{N} [x_1 \dots x_N] \cdot \begin{bmatrix} x_1^T \\ \vdots \\ x_N^T \end{bmatrix} = \frac{1}{N} X^T X$$

$N \times (p+1)$

Remark: Since $X = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_p \end{bmatrix}$
 $x_i = \begin{bmatrix} 1 \\ x_{i1} \\ \vdots \\ x_{ip} \end{bmatrix} \in \mathbb{R}^{p+1}$
 m
 \mathbb{R}^{p+1}

$$\bullet \mathbb{E}[XY] \approx \frac{1}{N} \sum_{i=1}^N x_i y_i = \frac{1}{N} \left(\int_{x_1} y_1 + \dots + \int_{x_N} y_N \right) = \frac{1}{N} \underbrace{\begin{bmatrix} x_1 & \dots & x_N \end{bmatrix}}_{X^T} \underbrace{\begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}}_y = \frac{1}{N} X^T y$$

Thus, β can be estimated:

$$\beta \approx \hat{\beta} = \left(\frac{1}{N} X^T X \right)^{-1} \cdot \frac{1}{N} X^T y = (X^T X)^{-1} X^T y \Rightarrow \boxed{\hat{\beta} = (X^T X)^{-1} X^T y} \quad (2)$$

So, if we assume that the relationship between the output Y and input X is approximately linear, $Y \approx f(X) = X^T \beta$, then we can estimate β using (2) and use $\hat{Y} = \hat{f}(X) = X^T \hat{\beta}$ for prediction. This is the linear regression method.
the best we can do under linear assumption.

If the linearity assumption is wrong, $Y \neq X^T \beta$, then the linear regression prediction $\hat{Y} = X^T \hat{\beta} \neq Y$ will obviously be not accurate.

Suppose that the assumption holds approximately: $Y = X^T \beta + \epsilon$, where $\epsilon \sim \mathcal{N}(0, \sigma^2)$

Gaussian noise

Q: How accurate the linear regression in this case?

Does the accuracy suffer from the curse of dimensionality?

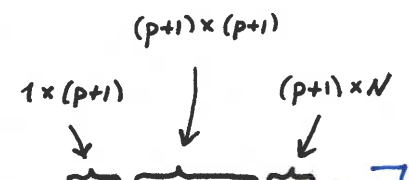
To answer these questions, let's derive the expected MSE of the method.

Let $X \in \mathbb{R}^{p+1}$ be a fixed input. The MSE at this X is

$$MSE(X) = \mathbb{E}[(Y - \hat{Y})^2] = \mathbb{E}[(X^T \beta + \epsilon - X^T \hat{\beta})^2] \leftarrow \mathbb{E} \text{ wrt } \epsilon \text{ and training data } \{X, y\}$$

$$\bullet \hat{\beta} = (X^T X)^{-1} X^T y = (X^T X)^{-1} X^T (X \beta + e) = \beta + (X^T X)^{-1} X^T e$$

$$\bullet y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} x_1^T \beta + \epsilon_1 \\ \vdots \\ x_N^T \beta + \epsilon_N \end{bmatrix} \Rightarrow y = \begin{bmatrix} x_1^T \\ \vdots \\ x_N^T \end{bmatrix} \beta + \underbrace{\begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_N \end{bmatrix}}_e = X \beta + e$$



$$\Rightarrow MSE(X) = \mathbb{E}[(X^T \beta + \epsilon - X^T \beta - X^T (X^T X)^{-1} X^T e)^2] = \mathbb{E}[(\epsilon - X^T (X^T X)^{-1} X^T e)^2]$$

$$= \mathbb{E}[\epsilon^2 - 2 \underbrace{\epsilon \cdot X^T (X^T X)^{-1} X^T e}_{\text{independent, } \mathbb{E}[\epsilon] = 0} + (X^T (X^T X)^{-1} X^T e)^2] = \sigma^2 + \underbrace{\mathbb{E}[(X^T (X^T X)^{-1} X^T e)^2]}_{\text{let's compute this expectation } \mathbb{E}}$$

$$\mathbb{E} = \mathbb{E}[X^T (X^T X)^{-1} X^T e \cdot \underbrace{X^T (X^T X)^{-1} X^T e}_{\text{scalar, we can transpose it}}] = \mathbb{E}[X^T (X^T X)^{-1} X^T e e^T X (X^T X)^{-1} X] =$$

$$= X^T \mathbb{E}[\mathbb{E}[(X^T X)^{-1} X^T e e^T X (X^T X)^{-1} | X]] X$$

$$= X^T \mathbb{E}[(X^T X)^{-1} X^T \underbrace{\mathbb{E}[e e^T | X]}_{= \mathbb{E}[e e^T] = \sigma^2 I_N} X (X^T X)^{-1}] X = X^T \mathbb{E}[\sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1}] X$$

\uparrow since e and X are independent
 $\mathbb{E}[e e^T] = \sigma^2 I_N$, where I_N is $N \times N$ identity matrix.

So, $E = \sigma^2 X^T E[(X^T X)^{-1}] X$

$\Rightarrow \boxed{MSE(X) = \sigma^2 + \sigma^2 X^T E[(X^T X)^{-1}] X}$

At the bottom of page 9, we saw that $\frac{1}{N} X^T X \approx E[XX^T]$.

Remark: By the law of large numbers $\frac{1}{N} X^T X \rightarrow E[XX^T]$ as $N \rightarrow \infty$.

Therefore, $(X^T X)^{-1} \approx \frac{1}{N} (E[XX^T])^{-1}$

Remark: If $E[X] = 0$, then $Cov(X) \stackrel{def}{=} E[(X - E[X])(X - E[X])^T] = E[XX^T]$

And, therefore, $(X^T X)^{-1} \approx \frac{1}{N} Cov(X)^{-1}$.

$\Rightarrow \boxed{MSE(X) \approx \sigma^2 + \sigma^2 X^T (E[XX^T])^{-1} X / N}$

Now we can find the expected MSE:

$E[MSE(X)] \approx \sigma^2 + \frac{\sigma^2}{N} E \left[\overset{\text{scalar}}{\underbrace{X^T E[XX^T]^{-1} X}_{\substack{(p+1) \times (p+1) \\ 1 \times (p+1) \quad (p+1) \times 1}}} \right] = \sigma^2 + \frac{\sigma^2}{N} E \left[\text{tr} \left(\underbrace{X^T E[XX^T]^{-1} X}_{\substack{\text{trace operator} \\ \text{tr } A = \sum_{i=1}^n A_{ii}}} \right) \right]$

$= \sigma^2 + \frac{\sigma^2}{N} E \left[\text{tr} (XX^T \cdot E[XX^T]^{-1}) \right]$

E and tr commute: $E[\text{tr}(X)] = \text{tr } E[X]$

$= \sigma^2 + \frac{\sigma^2}{N} \text{tr} \left(E[XX^T \cdot E[XX^T]^{-1}] \right)$

$= \sigma^2 + \frac{\sigma^2}{N} \text{tr} \left(E[XX^T] \cdot E[XX^T]^{-1} \right) = \sigma^2 + \frac{\sigma^2}{N} \text{tr } I_{p+1} = \sigma^2 + \frac{\sigma^2}{N} (p+1)$

Thus, we have: $\boxed{E[MSE(X)] \approx \sigma^2 + \frac{\sigma^2}{N} p}$

Under trace operator, matrix multiplication is a commutative operation:
 $\text{tr}(AB) = \text{tr}(BA)$
 Although $AB \neq BA$

As intuitively expected, prediction accuracy \uparrow as the sample size $N \uparrow$

Important observation:

We can suppress the curse of dimensionality by making σ^2/N small enough.

\downarrow as the dimension $p \uparrow$
 \downarrow as "noise strength" $\sigma \uparrow$

If σ^2/N is small, then the growth in the prediction error associated with the increase of p is negligible. (but this is only if $Y \approx X^T \beta$)