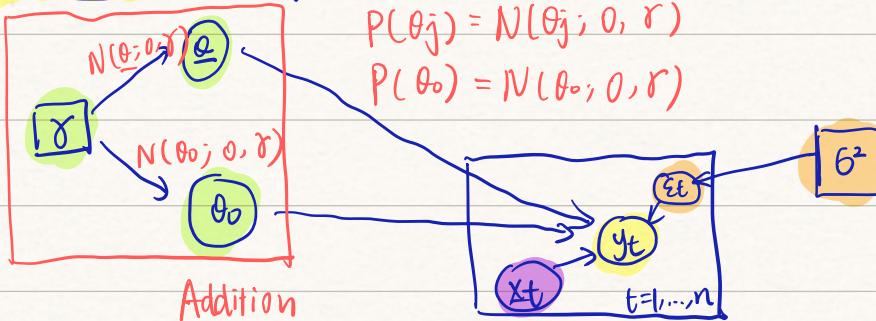


Graphical Model

(Bayesian Network)



$$P(\theta_j) = N(\theta_j; 0, r)$$

$$P(\theta_0) = N(\theta_0; 0, r)$$

- {
- Data
- Parameter
- Observation

$$\text{Generated Way} \Rightarrow y_t = \underline{\theta}^T x_t + \theta_0 + \varepsilon_t \quad \varepsilon_t \sim N(0, \sigma^2)$$

$$\begin{matrix} & r \\ \swarrow & \downarrow \\ r \end{matrix}$$

Take $(\underline{\theta}, \theta_0)$ as random variables, we have PENALIZED log-like likelihood.

$$\ell'(\underline{\theta}, \theta_0, \sigma^2, r) = \text{const} - \frac{n}{2} \log(\sigma^2) - \underbrace{\frac{1}{2\sigma^2} \sum_{t=1}^n (y_t - \underline{\theta}^T x_t - \theta_0)^2}_{\text{Original log-likelihood}}$$

$$- \underbrace{\frac{1}{2r} (\theta_0^2 + \sum_{j=1}^n \theta_j^2)}_{\text{Penalty term (on } \theta_0 \& \underline{\theta})} - \frac{d+1}{2} \log(r).$$

Penalty term (on θ_0 & $\underline{\theta}$)

$$\text{let } \frac{\partial \ell'}{\partial \theta} = 0 \Rightarrow \widehat{\underline{\theta}} = \begin{bmatrix} \widehat{\underline{\theta}} \\ \widehat{\theta}_0 \end{bmatrix} = (X^T X + \lambda I)^{-1} X^T \underline{y}.$$

$$\widehat{\theta} = [\underline{\theta}^T \theta_0]^T$$

Solution.

Model: (Ground true model)

$$y_t = (\theta^*)^T x_t + \theta_0^* + \varepsilon_t$$

$$\Leftrightarrow y = X \begin{bmatrix} \theta^* \\ \theta_0^* \end{bmatrix} + \varepsilon \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad X = \begin{bmatrix} x_1^T & 1 \\ \vdots & \vdots \\ x_n^T & 1 \end{bmatrix} \quad \varepsilon = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

{ Bias
Corvariance } of $\widehat{\theta}$ & θ^*
of ESTIMATOR & Ground True Parameter

$$\textcircled{1} \text{ Bias} = \mathbb{E} \left[\left[\frac{\hat{\theta}}{\theta_0} \right] - \left[\frac{\theta^*}{\theta_0^*} \right] \mid X \right]$$

$$\begin{aligned} \mathbb{E} \left[\left[\frac{\hat{\theta}}{\theta_0} \right] \mid X \right] &= (X^T X + \lambda I)^{-1} X^T \mathbb{E}[y \mid X] \\ &= (X^T X + \lambda I)^{-1} X^T \mathbb{E} \left[\left(X \left[\frac{\theta^*}{\theta_0^*} \right] + \varepsilon \right) \mid X \right] \\ &= (X^T X + \lambda I)^{-1} X^T X \left[\frac{\theta^*}{\theta_0^*} \right] \\ &= (X^T X + \lambda I)^{-1} (X^T X + \lambda I - \lambda I) \left[\frac{\theta^*}{\theta_0^*} \right] \\ &= \left[\frac{\theta^*}{\theta_0^*} \right] - \lambda (X^T X + \lambda I)^{-1} \left[\frac{\theta^*}{\theta_0^*} \right] = (I - \lambda (X^T X + \lambda I)^{-1}) \left[\frac{\theta^*}{\theta_0^*} \right] \end{aligned}$$

Therefore, $\mathbb{E} \left[\left[\frac{\hat{\theta}}{\theta_0} \right] - \left[\frac{\theta^*}{\theta_0^*} \right] \mid X \right] = -\lambda (X^T X + \lambda I)^{-1} \left[\frac{\theta^*}{\theta_0^*} \right]$

$$\Leftrightarrow \text{Bias} = -\lambda (X^T X + \lambda I)^{-1} \left[\frac{\theta^*}{\theta_0^*} \right]$$

Rmk: i) If $\lambda > 0$, there is a NON-ZERO Bias!

$$\text{ii) } \mathbb{E} \left[\left[\frac{\hat{\theta}}{\theta_0} \right] \mid X \right] := A \left[\frac{\theta^*}{\theta_0^*} \right] \quad A = \underbrace{I - \lambda (X^T X + \lambda I)^{-1}}$$

\textcircled{1} Notice that $0 \leq A \leq I$.

$\rightarrow A$ is positive semi-definite & all eigenvalue ≤ 1
that is to say $0 \leq \lambda_i(A) \leq 1 \quad i=1, 2, \dots, d+1$

\textcircled{2} Notice that $\mathbb{E} \left[\left[\frac{\hat{\theta}}{\theta_0} \right] \mid X \right] = A \left[\frac{\theta^*}{\theta_0^*} \right]$

$\left[\frac{\hat{\theta}}{\theta_0} \right]$ is the Regularized solution



$\left\{ \begin{array}{l} \text{tends to shrink the parameters!} \\ \text{bringing them close to 0!} \end{array} \right.$

Intuition: Regularization makes $(\hat{\theta}, \hat{\theta}_0)$ smaller!

$$(\hat{\theta})^T (\hat{\theta}) - (X^T X + \lambda I)^{-1} X^T u = (X^T X + \lambda I)^{-1} X^T X \left[\frac{\theta^*}{\theta_0^*} \right] + \varepsilon$$

② Covariance

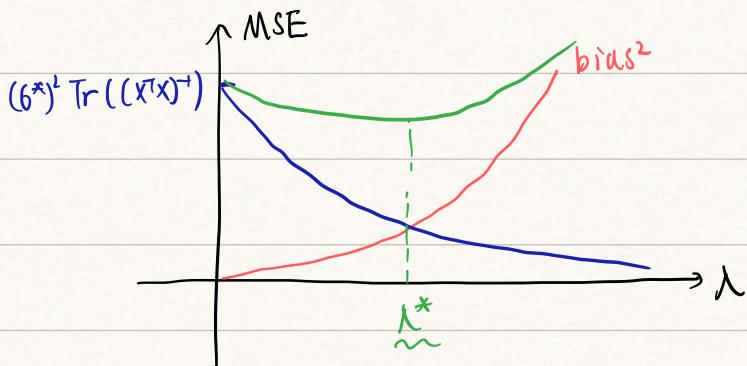
$$\text{Cov} \left(\begin{bmatrix} \hat{\theta} \\ \hat{\theta}_0 \end{bmatrix} \mid X \right) = \mathbb{E} \left[\left(\begin{bmatrix} \hat{\theta} \\ \hat{\theta}_0 \end{bmatrix} - \mathbb{E} \left[\begin{bmatrix} \hat{\theta} \\ \hat{\theta}_0 \end{bmatrix} \right] \right) \left(\begin{bmatrix} \hat{\theta} \\ \hat{\theta}_0 \end{bmatrix} - \mathbb{E} \left[\begin{bmatrix} \hat{\theta} \\ \hat{\theta}_0 \end{bmatrix} \right] \right)^T \mid X \right]$$

$$= (6^*)^2 \left[(X^T X + \lambda I)^{-1} - \lambda (X^T X + \lambda I)^{-2} \right]$$

$$\textcircled{3} \text{ MSE} = \|\text{BIAS}\|^2 + \text{Variance}$$

$$= \|\text{BIAS}\|^2 + \text{Tr}(\text{Covariance})$$

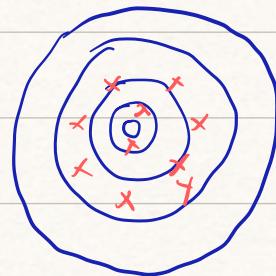
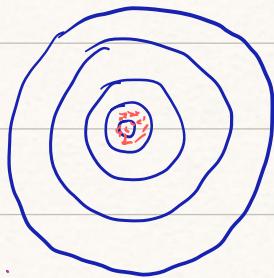
$$\left\{ \begin{array}{l} \text{Bias} = -\lambda (X^T X + \lambda I)^{-1} \left[\frac{\partial^X}{\partial \theta^*} \right] \\ \text{Cov} \left(\begin{bmatrix} \hat{\theta} \\ \hat{\theta}_0 \end{bmatrix} \mid X \right) = (6^*)^2 \left[(X^T X + \lambda I)^{-1} - \lambda (X^T X + \lambda I)^{-2} \right] \end{array} \right.$$



Bias-Variance

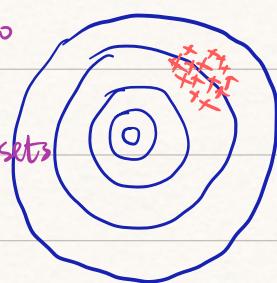
Low Bias & Low Variance

Low Bias, High Variance

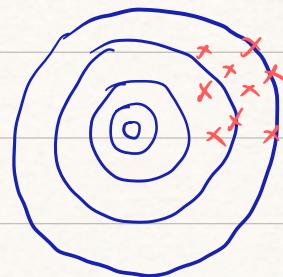


* Another example is
polynomial curve fitting

↓ # of data points
 { first order x
 3rd order ✓
 5th order x }
 n=8 High Bias, Low Variance
 N=20



High Bias, High Variance

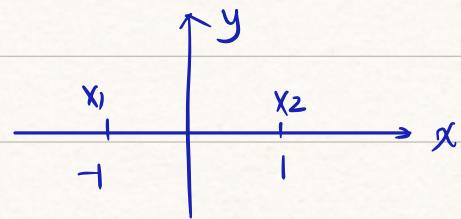


Example: $x_1 = -1, x_2 = +1$ $d = 1$

Design matrix $= \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} := X$

$$X^T X = 2I = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}$$

$$(\lambda I + X^T X) = (\lambda + 2) I = \begin{pmatrix} \lambda + 2 & 0 \\ 0 & \lambda + 2 \end{pmatrix}$$



$$\|\text{bias}\|^2 = \left(\frac{\lambda}{2+\lambda}\right)^2 \left[(\theta^*)^2 + (\theta_0^*)^2 \right]$$

$$\text{Variance} = (6^*)^2 \left(\frac{2}{2+\lambda} - \frac{2\lambda}{(2+\lambda)^2} \right)$$

$$\Rightarrow \text{MSE} = \frac{4(6^*)^2}{(2+\lambda)^2} + \frac{\lambda^2}{(2+\lambda)^2} \left((\theta^*)^2 + (\theta_0^*)^2 \right)$$

① If $\lambda=0$, $\text{MSE} = (6^*)^2$

② If $\lambda > 0$, and $(6^*)^2 > (\theta^*)^2 + (\theta_0^*)^2$

choose $\lambda = 2$. $\text{MSE}_{\lambda=2} < \frac{(6^*)^2}{2} < \text{MSE}_{\lambda=0}$.

↓
USE Regularization to Reduce MSE!

Active Learning : USE MSE to select future inputs!

① Given $\begin{cases} \underset{\downarrow}{x_1}, \underset{\downarrow}{x_2}, \dots, \underset{\downarrow}{x_n} \\ \underset{\downarrow}{y_1}, \underset{\downarrow}{y_2}, \dots, \underset{\downarrow}{y_n} \end{cases}$ choose x_{n+1} so that

MSE on $\{(x_t, y_t)\}_{t=1}^{n+1}$ is small.

(linear model)

→ Assume $y = \langle x, \theta^* \rangle + \theta_0^* + \varepsilon$ $\varepsilon \sim N(0, (6^*)^2)$

→ Assume X is non-random, so only source of randomness is y (through ε)

$$MSE = (\hat{b}^*)^2 \operatorname{Tr}((X^T X)^{-1}) \quad X = \begin{bmatrix} X_1^T & | & \\ \vdots & | & \\ X_n^T & | & \end{bmatrix} \in \mathbb{R}^{n \times (d+1)}$$

Criterion: Try to select X_{n+1} , which is appended to X (in the last row)

$$\tilde{X} = \begin{pmatrix} X & | & \\ \hline X_{n+1}^T & | & 1 \end{pmatrix} \in \mathbb{R}^{(n+1) \times (d+1)}$$

Calculate:

$$\tilde{X}^T \tilde{X} = \begin{bmatrix} X & | & \\ \hline X_{n+1}^T & | & 1 \end{bmatrix}^T \begin{bmatrix} X & | & \\ \hline X_{n+1}^T & | & 1 \end{bmatrix} = X^T X + \begin{bmatrix} X_{n+1} \\ 1 \end{bmatrix} \begin{bmatrix} X_{n+1} \\ 1 \end{bmatrix}^T$$

$$:= A^{-1} + VV^T$$

inverse

Notation: $\begin{cases} A = (X^T X)^{-1} \rightarrow \text{old Gram Matrix} \\ V = \begin{bmatrix} X_{n+1} \\ 1 \end{bmatrix} \rightarrow \text{new data vector} \end{cases}$

$$\underbrace{\text{New MSE}}_{\text{MSE}'} = (\hat{b}^*)^2 \operatorname{Tr}(\underbrace{(A^{-1} + VV^T)^{-1}}_{\text{Matrix}})$$

Goal: Find a "Valid" $\underline{x}_{n+1} \in \mathbb{R}^d$ (i.e., $V \in \mathbb{R}^{d+1}$), such that MSE' is as small as possible.

$$\text{Fact: } (A^{-1} + VV^T) = A - \frac{AVV^TA}{1 + V^TAV} \quad \begin{array}{l} \text{Matrix} \\ \text{Scalar} \end{array}$$

$$\text{MSE}' \propto \operatorname{Tr}((A^{-1} + VV^T)^{-1})$$

$$\text{New MSE} = \text{Tr}(A) - \frac{1}{1 + V^T A V} \text{Tr}(V^T A A V)$$

$$= \underbrace{\text{Tr}(A)}_{\text{Old MSE}} - \underbrace{\frac{1}{1 + V^T A V} \text{Tr}(V^T A A V)}_{\text{positive!}}$$

nonnegative!
Fact: $V^T A A V \geq 0$ for all V
 $\Leftrightarrow \|AV\|^2 \geq 0$

Fact: $V^T A V \geq 0$ for all V

PF: $A = (X^T X)^{-1}$

$X^T X$ is positive semi-definite

↓
positive definite (inverse)

$\Rightarrow (X^T X)^{-1}$ positive definite

$\Leftrightarrow A$ positive definite

Rmk: for $\forall V \in \mathbb{R}^{d+1}$, MSE' \leq MSE!

Goal: Find X that maximizes the Reduction (Of MSE)

$$\frac{V^T A A V}{1 + V^T A V}$$

$V = \begin{pmatrix} X \\ 1 \end{pmatrix}$

$$\text{MSE}_{\text{original}} = (\hat{b}^*)^2 \text{Tr}((X^T X)^{-1}) = (\hat{b}^*)^2 \text{Tr}(A)$$

$$\text{MSE}_{n+1} = (\hat{b}^*)^2 \left[\text{Tr}(A) - \frac{V^T A A V}{1 + V^T A V} \right]$$

$$= \text{MSE}_{\text{origin}} - (\hat{b}^*)^2 \cdot \text{Reduction}$$

Maximize this!

Note:

No matter how smart you are.

Claim: $\frac{\underline{V}^T A A \underline{V}}{1 + \underline{V}^T A \underline{V}} \leq \lambda_{\max}(A)$ the reduction has an
Upper Bound!

and the upper bound in (*) is attained when $\underline{V} = t \underline{u}_1$ as $t \rightarrow \infty$.

\underline{u}_1 is the eigenvector corresponding to $\lambda_{\max}(A)$.

(i.e., the top eigenvalue of A).

Rmk: A is symmetric P.S.D.

$$(A \in \mathbb{R}^{(d+1) \times (d+1)})$$

$$\boxed{A} \left\{ \begin{array}{l} \lambda_1 \rightarrow \underline{u}_1 \\ \vdots \\ \lambda_{d+1} \rightarrow \underline{u}_{d+1} \\ \{ \underline{u}_i \} \text{ Orthonormal} \end{array} \right.$$

Rmk: $\lim_{t \rightarrow \infty} \frac{\underline{V}^T A A \underline{V}}{1 + \underline{V}^T A \underline{V}} \Big|_{\underline{V} = t \underline{u}_1} = \lambda_{\max}(A)$

In conclusion, to max the reduction, choose \underline{V} to be 'aligned' with \underline{u}_1 and have the largest possible magnitude.

For example, constrain $\|\underline{V}\| \leq C$ for some $0 < C < \infty$.

we choose $\underbrace{t = C}_{\text{①}} \quad \underbrace{\underline{u}_1 \rightarrow \text{the normalized eigenvector (top)}}_{\text{②}}$

$$\boxed{\underline{V} = t \underline{u}_1}$$

Pf. $\{\underline{u}_i\}_{i=1}^{d+1}$ forms as basis of A . (orthogonal eigenvectors)

Every $\underline{V} \in \mathbb{R}^{d+1}$ can be written as $\underline{V} = \sum_{i=1}^{d+1} \alpha_i \underline{u}_i$

$$A \underline{V} = \sum_i \alpha_i A \underline{u}_i = \sum_i \alpha_i \lambda_i \underline{u}_i$$

$$\textcircled{1} \quad \underline{V}^T A \underline{V} = \left(\sum_i \alpha_i \underline{u}_i \right)^T \left(\sum_j \alpha_j \lambda_j \underline{u}_j \right)$$

$$= \sum_i \alpha_i^2 \lambda_i^2 \quad \left\{ \begin{array}{l} \underline{u}_i^T \underline{u}_j = 0, \quad i \neq j \\ \underline{u}_i^T \underline{u}_i = 1, \quad i = j \end{array} \right.$$

Orthogonal

$$\textcircled{2} \quad \underline{V}^T A A \underline{V} = \sum_i \alpha_i^2 \lambda_i^2$$

$$\text{then } \frac{\sqrt{V^T A A V}}{1 + \sqrt{V^T A V}} \leq \frac{\sqrt{V^T A V}}{\sqrt{V^T V}} \leq \frac{\sum_i \lambda_i^2 \lambda_i}{\sum_i \lambda_i^2} \leq \frac{\lambda_1 \sum_i \lambda_i^2 \lambda_i}{\sum_i \lambda_i^2 \lambda_i} = \lambda_1.$$

#

$$\text{Example: } x_1 = -1 \quad x_2 = +1$$

$$y = \theta_0 x + \theta_1 + \varepsilon \quad x \in [-1, 1].$$

$$V = \begin{bmatrix} * \\ 1 \end{bmatrix} \quad X = \begin{bmatrix} -1 & 1 \\ 1 & 1 \end{bmatrix} \quad X^T X = 2I.$$

$$A = (X^T X)^{-1} = \frac{1}{2} I$$

$$\textcircled{1} \quad V^T A V = \begin{pmatrix} * \\ 1 \end{pmatrix}^T \frac{1}{2} I \begin{pmatrix} * \\ 1 \end{pmatrix} = \frac{x^2 + 1}{2}$$

$$\textcircled{2} \quad V^T A A V = \frac{x^2 + 1}{4}$$

Want to maximize the Reduction in MSE

$$\text{Reduction} = \frac{\sqrt{V^T A A V}}{1 + \sqrt{V^T A V}} = \frac{\frac{1+x^2}{4}}{1 + \frac{1+x^2}{2}} = \frac{1}{2} \frac{z}{1+z} \quad (z = \frac{1+x^2}{2})$$

$$\frac{z}{1+z} = 1 - \frac{1}{1+z} \quad \uparrow z$$

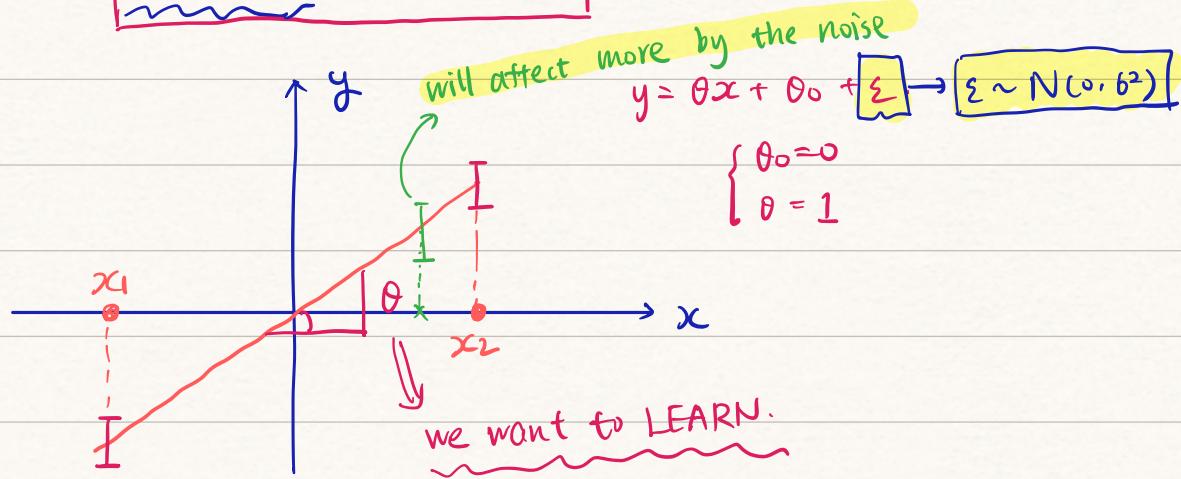
\Rightarrow choose z as big as possible!

\Rightarrow choose $|x|$ as large as possible!

\rightarrow Conclusion: choose x (next point) to be ± 1

Rmk: For linear models inputs should be as far apart as

possible to amortize the noise

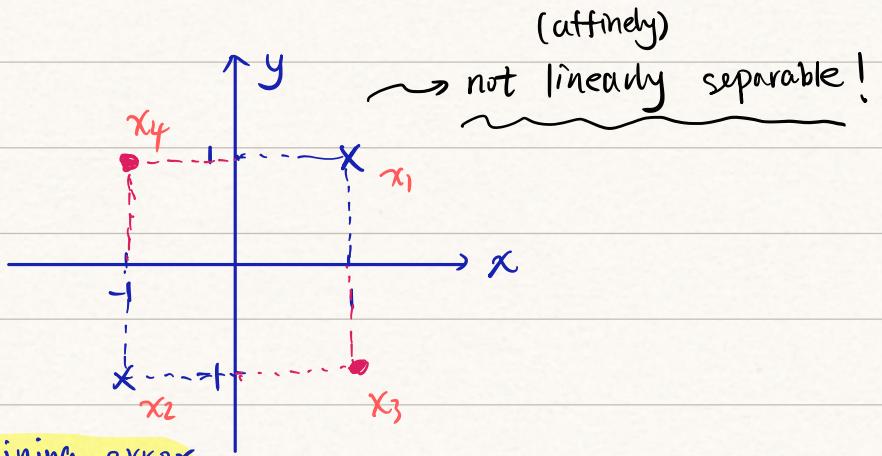


Introduction to KERNEL

XOR dataset

* there exists NO AFFINE

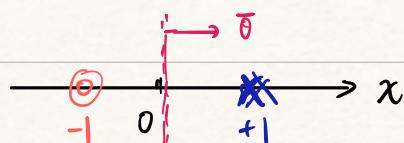
classifiers that attain 0 training error



↓
Construct a New Dataset (by multiplying the 2 coordinations together)

$$X'_1 = X_{11} X_{12} = 1 \quad X'_2 = X_{21} X_{22} = 1 \quad \rightarrow \text{New Positive Label}$$

$$X'_3 = X_{31} X_{32} = -1 \quad X'_4 = X_{41} X_{42} = -1 \quad \rightarrow \text{New Negative Label}$$



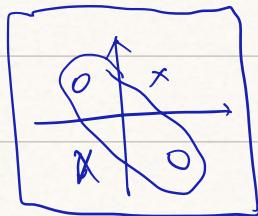
A decision boundary at $x=0$ attains 0 training error!

1. Consider the linear model : $y = \theta_0 + \theta_1 x + \epsilon$ $\epsilon \sim N(0, \sigma^2)$

$$\textcircled{1} \quad \phi: x \in \mathbb{R} \rightarrow \phi(x) = \begin{pmatrix} \sqrt{2}x \\ x^2 \end{pmatrix} \rightsquigarrow \text{quadratic kernel}$$

$$\textcircled{2} \quad \phi: x \in \mathbb{R} \rightarrow \phi(x) = \begin{pmatrix} \sqrt{3}x \\ \sqrt{3}x^2 \\ x^3 \end{pmatrix} \rightsquigarrow \text{cubic kernel.}$$

This RESULT IN a polynomial regression model !



$$y = \langle \underline{\theta}, \underline{\phi}(x) \rangle + \theta_0 + \epsilon \Rightarrow \text{Actually, contains all polynomials whose order } \leq 3 \text{ (cubic kernel)}$$

2. Polynomial expansion also works in Higher dimension

$$\hookrightarrow \text{i.e., } \underline{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \Leftrightarrow \text{e.g. XOR dataset}$$

$$\underline{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \in \mathbb{R}^2 \xrightarrow[\substack{\text{quadratic} \\ \text{kernel}}]{\phi} \underline{\psi}(\underline{x}) = \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ x_1 x_2 \\ x_1^2 \\ x_2^2 \end{pmatrix} \in \mathbb{R}^6 \rightarrow \text{govern by} \begin{cases} d: \text{input dimension} \\ p: \text{kernel function} \end{cases}$$

Rmk: $\underline{\psi}(\underline{x})$ contains all monomials whose orders ≤ 2

(as for quadratic kernel)

defn. monomials order $\leq p$.

$$\Leftrightarrow x_1^a x_2^b \text{ such that } a+b \leq p$$

Note: Dimension of $\underline{\psi}(\underline{x})$ grows rapidly with (d, p)

$$\text{Claim: } \dim (\underline{\psi}(\underline{x})) = \binom{p}{d+1}$$

Pf: Suppose that $\underline{x} \in \mathbb{R}^d$, namely $\begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix}$

and $\underbrace{x_1^{N_1} \cdots x_d^{N_d}}$ find the nonnegative integer

solution $(\lambda_1, \dots, \lambda_d)$ such that $\sum_{i=1}^d \lambda_i \leq p$

$\Leftrightarrow \sum_{i=1}^d \lambda_i = k + d \rightarrow (k) \quad \lambda_i \geq 1, \quad \lambda_i \in \mathbb{Z}. \quad k = 0, 1, 2, \dots, p.$

The # of feasible solution for (P) is $\binom{d-1}{k+d-1}$

$$\begin{aligned} \Leftrightarrow \sum_{k=0}^p \binom{d-1}{k+d-1} &= \binom{d-1}{d-1} + \binom{d-1}{d} + \cdots + \binom{d-1}{d+p-1} \\ &= \binom{d-1}{d} + \binom{d-1}{d} + \cdots + \binom{d-1}{d+p-1} \\ &= \binom{d}{d+1} + \binom{d-1}{d+1} + \cdots + \binom{d-1}{d+p-1} \\ &= \binom{d}{d+2} + \cdots + \binom{d-1}{d+p-1} \\ &= \cdots = \binom{d}{d+p} \quad \square \end{aligned}$$

Q: Should we first form the LONG FEATURE VECTOR $\underline{\phi}(\underline{x})$ [With dimension $\binom{d}{p+d}$] , then take inner product?

$$\langle \underline{\phi}(\underline{x}), \underline{\phi}(\underline{x}') \rangle$$

Consider: $\underline{x}, \underline{x}' \in \mathbb{R}$, & $\underline{\phi}: x \in \mathbb{R} \rightarrow (1, \sqrt{3}x, \sqrt{3}x^2, x^3)^T$

$$\text{then } \langle \underline{\phi(x)}, \underline{\phi(x')} \rangle = \begin{pmatrix} \sqrt{3}x \\ \sqrt{3}x^2 \\ x^3 \end{pmatrix}^\top \begin{pmatrix} \sqrt{3}x' \\ \sqrt{3}x'^2 \\ x'^3 \end{pmatrix}$$

$$= 1 + 3xx' + 3(x'x)^2 + (x'x)^3$$

$$= (1+x'x)^3 !$$

Observation : if we denote $\underline{x}'\underline{x} = \langle \underline{x}, \underline{x}' \rangle_{\mathbb{R}} = z$

$$\langle \phi(x), \phi(x') \rangle_{\mathbb{R}^4} = (1+z)^3$$

Therefore. Can we evaluate the inner product $\langle \underline{x}, \underline{x}' \rangle_{\mathbb{R}}$

first, then passing it to the inner product $\langle \underline{\phi(x)}, \underline{\phi(x')} \rangle_{\mathbb{R}^d}$

through a simple function (like $f(z) = (1+z)^3$)

if we : ① $\langle \underline{x}, \underline{x}' \rangle_{\mathbb{R}^d}$

② $\langle \underline{\phi(x)}, \underline{\phi(x')} \rangle_{\mathbb{R}^d} = f(\langle \underline{x}, \underline{x}' \rangle)$

difficult to calculate directly !

$O(d^3)$

Naive Computation

Smart Computation !

easy to calculate ! (only takes time $O(d)$)