

# REGULARIZATION

## Nonlinear Feature Maps

One way to create nonlinear estimators or classifiers is to first transform the data via a nonlinear feature map

$$\Phi: \mathbb{R}^d \rightarrow \mathcal{H}$$

and apply a linear method to the transformed data  $\Phi(x_1), \dots, \Phi(x_n)$

Example 1 |  $y_i = f(x_i) + \varepsilon_i, \quad i=1, \dots, n$

$$f(x) = \sum_{j=0}^p \beta^{(j)} \cdot x^j \quad (\text{degree } p \text{ polynomial})$$

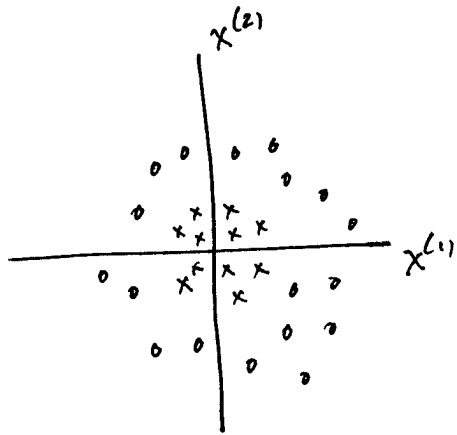
To estimate  $f$ , we may apply least-squares regression to the transformed data  $\Phi(x_i)$ ,

where  $\Phi(x) = [1 \ x \ x^2 \ \dots \ x^p]^T$

$$\Rightarrow \hat{\beta} = (A^T A)^{-1} A^T \underline{y}, \quad A =$$

Ⓐ

## Example 2



$$x \mapsto \Phi(x) = \begin{bmatrix} 1 \\ x^{(1)} \\ x^{(2)} \\ x^{(1)} \cdot x^{(2)} \\ (x^{(1)})^2 \\ (x^{(2)})^2 \end{bmatrix}$$

Then the data are linearly separable in the new feature space. They are correctly classified by

$$\text{sign} \{ w^T \Phi(x) \}$$

where

(B)

$$w =$$

In many applications, we don't know exactly how to design  $\Phi(x)$ , we just know that some nonlinear features are probably important.

In such situations, it is common to include a large number of nonlinear features, in hopes that some of them are relevant.

Unfortunately, this practice can lead to \_\_\_\_\_

© \_\_\_\_\_ problems.

Example 1, revisited In least squares polynomial regression, we have

$$\hat{\beta} = (A^T A)^{-1} A^T \underline{y}$$

where

$$A = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^p \\ 1 & x_2 & x_2^2 & \dots & x_2^p \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^p \end{bmatrix}$$

As  $p$  increases, the smallest eigenvalue of  $A^T A$  gets very, very small, so that matrix

① inversion is numerically \_\_\_\_\_.

Alternatively, the least-squares criterion becomes very \_\_\_\_\_ near the minimizer.

Essentially, when there are many features, it is possible that a huge coefficient on one feature could be cancelled out by other features.

To remedy the situation, we can incorporate a regularization term into design criteria that will keep coefficients \_\_\_\_\_.

Owing to computational convenience, the most common kind of regularization is \_\_\_\_\_.

We'll examine two cases in detail.

## Ridge Regression

Given  $y_i = f(x_i) + \epsilon_i$ , where  $f(x_i) = \beta^T x_i + \beta_0$ .

Instead of minimizing the sum of squared errors, in ridge regression, we minimize

$$\sum_{i=1}^n (y_i - \beta^T x_i - \beta_0)^2 + \lambda \|\beta\|^2$$

where  $\lambda > 0$  is a tuning parameter.

Note:  $\beta_0$  is not penalized so that our solution is independent of where the origin is located.

Let's derive the solution. First, let's eliminate  $\beta_0$

$$\frac{\partial}{\partial \beta_0} = -2 \sum_{i=1}^n (y_i - \beta^T x_i - \beta_0) = 0$$

$$\Rightarrow \hat{\beta}_0 = \frac{1}{n} \sum_i y_i - \hat{\beta}^T x_i$$

①

=

Thus, we are left to minimize

$$\sum_{i=1}^n (y_i - \bar{y} - \beta^T (x_i - \bar{x}))^2 + \lambda \beta^T \beta$$

wrt  $\beta$ . For convenience, assume  $\bar{y} = 0$ ,  $\bar{x} = 0$ .

The criterion may be written

$$(\underline{y} - A\underline{\beta})^T (\underline{y} - A\underline{\beta}) + \lambda \underline{\beta}^T \underline{\beta}, \quad A = \begin{bmatrix} x_1^{(1)} & \dots & x_1^{(n)} \\ \vdots & & \vdots \\ x_n^{(1)} & \dots & x_n^{(n)} \end{bmatrix}$$

$$= \underline{y}^T \underline{y} + \underline{\beta}^T A^T A \underline{\beta} - 2 \underline{\beta}^T A^T \underline{y} + \lambda \underline{\beta}^T \underline{\beta}$$

$$= \underline{\beta}^T [A^T A + \lambda I] \underline{\beta} - 2 \underline{\beta}^T A^T \underline{y} + \underline{y}^T \underline{y}$$

$$\frac{\partial}{\partial \underline{\beta}} = 0 \implies (A^T A + \lambda I) \underline{\beta} = A^T \underline{y}$$

Ⓕ  $\implies$

Observations •  $\lambda = 0$  recovers least-squares linear regr.

- $\lambda I$  increases the eigenvalues of  $A^T A$  by  $\lambda$ , so that  $A^T A + \lambda I$  is not ill-conditioned.

## Soft Margin Hyperplane

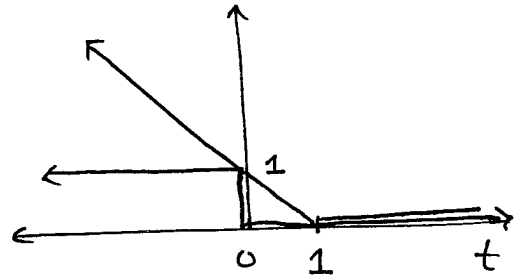
The training error of a linear classifier may be bounded as

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{y_i(\omega^\top x_i + b) < 0\}} \\ \leq \frac{1}{n} \sum \phi(y_i(\omega^\top x_i + b)) \end{aligned}$$

where  $\phi(t)$  is any upper bound on  $\mathbb{1}_{\{t < 0\}}$ .

Let's take

$$\begin{aligned} \phi(t) &= \max\{0, 1-t\} \\ &=: (1-t)_+ \end{aligned}$$



In addition, let's add a quadratic penalty to keep the coefficients small.

$$\Rightarrow \min_{w, b} \frac{1}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n (1 - y_i(\omega^\top x_i + b))_+$$

Compare this to the quadratic program for the optimal soft-margin hyperplane:

$$\min_{w, b, \xi_i} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

$$\text{s.t. } y_i (w^T x_i + b) \geq 1 - \xi_i, \quad i=1, \dots, n$$

$$\xi_i \geq 0, \quad i=1, \dots, n$$

Claim If  $C = \frac{1}{n\lambda}$ , these two optimization problems are solved by the same  $w, b$ .

⑥ Proof:



Key

A.

$$A = \begin{bmatrix} \Phi(x_1)^T \\ \Phi(x_2)^T \\ \vdots \\ \Phi(x_n)^T \end{bmatrix}$$

B.

$$w = \begin{bmatrix} -r^2 \\ 0 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$$

(circle of radius  $r$ )

C. ill-conditioned  
quadratic

D. unstable, flat, small,

E.  $\bar{y} = \hat{\beta}^T \bar{x}$

F.  $\hat{\beta} = (A^T A + \lambda I)^{-1} A^T y$

G. If  $\xi_i > 0$ , then  $y_i(w^T x_i + b) = 1 - \xi_i$ ,

If  $\xi_i = 0$ , then  $y_i(w^T x_i + b) \geq 1$

$$\Rightarrow \sum \xi_i = \sum (1 - y_i(w^T x_i + b))$$

at the global  
minimizer

if not, we could decrease  
the objective function without  
violating the constraints