

Linear Regression

" $y = wx + c$ " offic math or " $X_{aug} \cdot \hat{w} = y$ " matrix eqn

$$y = \begin{pmatrix} 1 & x \end{pmatrix} \begin{pmatrix} c \\ w \end{pmatrix}$$

X_{aug} w

e.g. $f(x_1, x_2) = w_0 + w_1 x_1 + w_2 x_2$

where $f(\cdot)$ predicts y using x_1 & x_2 .

Available data: X_{train} , y_{train} , X_{test}

what to do?

- ① Augment X : $X_{-a} \leftarrow \begin{pmatrix} 1 \\ \vdots \end{pmatrix} X_{train}$
- ② Find \hat{w} .
- ②A determine shape of X_{-a} : $m = \text{no. of rows (samples)}$
 $d = \text{no. of coln. (features)}$
- ②B Apply Left / Right Inverse to get linear weights for the linear model

$m > d$ (primal): $\hat{w} \leftarrow (X_{-a}^T X_{-a})^{-1} X_{-a}^T \cdot y_{train}$ (Left inverse)

$m < d$ (dual): $\hat{w} \leftarrow X_{-a}^T (X_{-a} X_{-a}^T)^{-1} \cdot y_{train}$ (Right inverse)

constant coeff

Note: $\hat{w} = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{pmatrix}$

w_i is the weight for x_i

$X = \begin{pmatrix} x_1, \dots, x_n \end{pmatrix}$

③ Apply linear model to X_{test} to predict y_{pred}

! need to augment X_{test} first!

$X_{-a} = X_{aug}$

a) $X_{test-aug} \leftarrow \begin{pmatrix} 1 \\ \vdots \end{pmatrix} X_{test}$

b) $y_{pred} \leftarrow X_{test-aug} \cdot \hat{w}$

Opt: 4 compare y_{pred} with a hidden y_{test} : compute $MSE(y_{pred}, y_{test})$ } test error.

Ridge Regression

What happens if $X^T X$ (resp. $X X^T$) is not invertible?

→ use $X^T X + \lambda I$ (resp. $X X^T + \lambda I$) instead with $\lambda > 0$ small. e.g. $\lambda = 0.0001$

→ this is called regularization.

→ two purposes: ① makes $X^T X + \lambda I$ invertible
② biases the cost fn, to prefer \hat{w} that is smaller

cost fn in primal form

original cost: $\min_w \|Xw - y\|^2 = \min_w (Xw - y)^T (Xw - y)$ $\|\hat{w}\|^2 = \hat{w}^T \hat{w}$ is small

regularized cost: $\min_w \|Xw - y\|^2 + \lambda \|w\|^2 = \min_w (Xw - y)^T (Xw - y) + \lambda w^T w$

How to get w that minimizes the cost? • differentiate to get soln → Left Inverse. is now

$(X^T X + \lambda I)^{-1} X^T$

cost fn in dual form (extra)

original cost: $\min_w \|Xw - y\|^2 = \min_w \sum_{i=1}^m (y_i - x^{(i)} w)^2$ where $X = \begin{pmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(m)} \end{pmatrix}$, $y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix}$

regularized cost: $\min_w \|Xw - y\|^2 + \lambda \|w\|^2 = \min_w \sum_{i=1}^m (y_i - x^{(i)} w)^2 + \lambda \sum_{j=1}^d w_j^2$, where $w = \begin{pmatrix} w_1 \\ \vdots \\ w_d \end{pmatrix}$

How to get (extra)

w that minimizes the cost? • need to apply Lagrangian duality. See KKT theory. This is why this is extra.

rough idea: get dual problem
solve dual problem.

Ridge Regression (what to do? SAME as Linear Regression but change the Left/Right Inverse)

- ① Augment X : $X_{\text{aug}} \leftarrow \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} X_{\text{train}}$
- ② determine shape of X_{aug} : $m = \text{no. of rows}$
 $d = \text{no. of cols.}$
- ③ Apply Left/Right Inverse to get linear weights for the linear model
 - $m > d$ (primal): $\hat{w} \leftarrow (X^T X + \lambda I)^{-1} X^T \cdot y_{\text{train}}$ (Left inverse)
 - $m < d$ (dual): $\hat{w} \leftarrow X^T (X X^T + \lambda I)^{-1} \cdot y_{\text{train}}$ (Right inverse)
- ④ Apply linear model to X_{test} to predict y_{pred}
 - ! need to augment X_{test} first!
 - $X := X_{\text{aug}}$

constant coeff

Note: $\hat{w} = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{pmatrix}$

w_0 is the weight for x_0

$X = \begin{pmatrix} x_1 & \dots & x_n \end{pmatrix}$

Polynomial Regression (examples)

→ non-linear model

→ e.g. $f(x) = w_0 + w_1 x + w_2 x^2 + \dots + w_k x^k$. This is a k^{th} order polynomial in one variable x .

given x_{train} , we want to

estimate weights w_0, w_1, \dots, w_k , which are the "linear weights for each monomial" $x^i \sim w_i$.

→ e.g. $f(x_1, x_2) = w_{0,0} + w_{1,0} x_1 + w_{0,1} x_2 + w_{2,0} x_1^2 + w_{1,1} x_1 x_2 + w_{0,2} x_2^2$

$x_{\text{train}} = (x_1, x_2)$

This is a 2nd-order polynomial model in two variables, x_1 & x_2 .

observe that $f(x_1, x_2) = \begin{pmatrix} 1 & x_1 & x_2 & x_1^2 & x_1 x_2 & x_2^2 \end{pmatrix}$

so we can express our polynomial model LINEARLY with respect to monomials

$\begin{pmatrix} w_{0,0} \\ w_{1,0} \\ w_{0,1} \\ w_{2,0} \\ w_{1,1} \\ w_{0,2} \end{pmatrix} = P_{\text{train}} \cdot \hat{w}$

Model: $P_{\text{train}} \hat{w} = y$,

where for each sample of X_{train} , e.g. $(x_1, x_2) = (2, 5)$

we replace with $P_{\text{train}} \leftarrow \begin{pmatrix} 1 & x_1 & x_2 & x_1^2 & x_1 x_2 & x_2^2 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{pmatrix}$
 $= \begin{pmatrix} 1 & 2 & 5 & 4 & 10 & 25 \end{pmatrix}$

! Notice that P_{train} is already augmented. the "1" correspond to constant coeff $w_{0,0}$.
& $x_1^i x_2^j$ correspond to $w_{i,j}$.

Polynomial Regression (Steps)

Data: $X_{\text{train}}, y_{\text{train}}, X_{\text{test}}$

what to do?

① : transform X_{train} to P_{train} row by row.

e.g. $X_{\text{train}} = \begin{pmatrix} x_{1,1} & x_{1,2} \\ x_{2,1} & x_{2,2} \\ x_{3,1} & x_{3,2} \end{pmatrix} \rightarrow \begin{pmatrix} 1 & x_{1,1} & x_{2,1} & x_{1,1}^2 & x_{1,1}x_{2,1} & x_{1,2}^2 \\ 1 & x_{2,1} & x_{2,1} & x_{2,1}^2 & x_{2,1}x_{2,2} & x_{2,2}^2 \\ 1 & x_{3,1} & x_{3,1} & x_{3,1}^2 & x_{3,1}x_{3,2} & x_{3,2}^2 \end{pmatrix} = P_{\text{train}}$

\uparrow constant \uparrow x_1 \uparrow x_2 \uparrow x_1^2 \uparrow x_1x_2 \uparrow x_2^2

②A Determine shape of P_{train}

$m = ?$ same as X .

$d = ?$ how many monomials?

M1: count one-by-one

M2: use formula

$$\binom{d+k}{k}$$

where we have d variables $\{x_1, x_2, \dots, x_d\}$ & k^{th} order polynomial model

②B Apply Left / Right Inverse to get linear weights for the polynomial model

$m > d$ (primal): $\hat{w} \leftarrow (P^T P + \lambda I)^{-1} P^T \cdot y_{\text{train}}$ (Left inverse)
 $m < d$ (dual): $\hat{w} \leftarrow P^T (P P^T + \lambda I)^{-1} \cdot y_{\text{train}}$ (Right inverse)

\hat{w} looks like $\begin{pmatrix} w_{0,0} \\ w_{1,0} \\ \vdots \\ w_{1,1} \end{pmatrix}$

④ Apply polynomial model to X_{test} to predict y_{pred}

! need to transform X_{test} to P_{test} first

⑥ P_{test} gotten similarly to Step ①

⑥ $y_{\text{pred}} \leftarrow P_{\text{test}} \cdot \hat{w}$! This is same as $y_{\text{pred}} = f(x_{\text{test}})$

Binary Classification "class = $\text{sgn}(Xw) = \text{sgn}(y)$ "

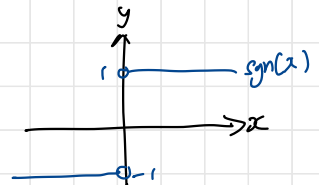
Data: $X_{\text{train}}, \text{class}_{\text{train}}, X_{\text{test}}$, where $\text{class}_{\text{train}} = \begin{pmatrix} \pm 1 \\ \vdots \\ \pm 1 \end{pmatrix} \in \{\pm 1\}^m$

what to do?

① Apply Linear or Polynomial Regression to $X_{\text{train}}, \text{class}_{\text{train}}, X_{\text{test}}$ to get $y_{\text{pred}} = X_{\text{test}} \cdot \hat{w}$ (or $y_{\text{pred}} = P_{\text{test}} \cdot \hat{w}$)

② Apply sigm-function to y_{pred}

\downarrow
 $\text{class}_{\text{pred}} = \text{sgn}(y_{\text{pred}})$
 set $y=0$ as threshold.



Multi-Category classification "class = $\arg\max_i Xw$ "

Data: X_{train} , class_train , X_{test} , where $\text{class_train} = \begin{pmatrix} 1 \\ 2 \\ 1 \\ 3 \end{pmatrix} \in \{1, 2, \dots, c\}^m$
! class_train has 'c' amount of classes.

What to do?

Step ①: Apply one-hot-encoding $y_{\text{train}} = \text{OHE}(\text{class_train})$.

idea: sample by sample, convert each output in class to a standard basis vector in \mathbb{R}^c

e.g. $c=3$

$$\text{class_train} = \begin{pmatrix} 1 \\ 2 \\ 1 \\ 3 \end{pmatrix} \xrightarrow{\text{OHE}} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} = y_{\text{train}}$$

! observe that y_{train} is a (m, c) array

Step ②: Perform Linear Regression OR Polynomial Regression on X_{train} , y_{train} , X_{test} .

so in linear reg. " $Xw = y$ "

w is a (d, c) -array

→ get y_{pred} .

e.g. $y_{\text{pred}} = \begin{pmatrix} 0.5 & 0.1 & 0.2 \\ 0.1 & 0.3 & 0.4 \\ 0.3 & 0.2 & 0.7 \\ 0.8 & 0.2 & 0.1 \end{pmatrix}$

Step ③: Get the class that is 'most probable'

e.g. $y_{\text{pred}} = \begin{pmatrix} 0.5 & 0.1 & 0.2 \\ 0.1 & 0.3 & 0.4 \\ 0.3 & 0.2 & 0.7 \\ 0.8 & 0.2 & 0.1 \end{pmatrix} \xrightarrow{\text{argmax}} \begin{pmatrix} 1 \\ 3 \\ 3 \\ 1 \end{pmatrix} = \text{class_pred}$

EE2211 Tutorial 6

(Ridge Regression in Dual Form)

Question 1:

Derive the solution for linear ridge regression in dual form (see Lecture 6 notes page 16).

(Polynomial Regression, 1D data)

Question 2:

Given the following data pairs for training

- $\{x = -10\} \rightarrow \{y = 5\}$
- $\{x = -8\} \rightarrow \{y = 5\}$
- $\{x = -3\} \rightarrow \{y = 4\}$
- $\{x = -1\} \rightarrow \{y = 3\}$
- $\{x = 2\} \rightarrow \{y = 2\}$
- $\{x = 8\} \rightarrow \{y = 2\}$

$$\begin{matrix} * & ? \\ x_0 = -10 & \\ x_1 = -8 & \end{matrix} \begin{pmatrix} 1 & -10 & (-10)^2 & (-10)^3 \\ 1 & -8 & (-8)^2 & (-8)^3 \end{pmatrix} \quad x = -10.$$

linear model
 $f(x) = w_0 + w_1 x + w_2 x^2$
 $x \mapsto \boxed{\quad} \rightarrow y$

How many inputs: 1
 shape of X : $(6, 1)$
 polynomial model: $f(x) = w_0 + w_1 x + w_2 x^2 + w_3 x^3$
 How to get P ? : *
 Is our polynomial regression primal or dual form?
 \rightarrow how many variables? P is of shape $(6, 4)$

$$\begin{matrix} m_p = 6 \\ d_p = 4 \end{matrix} \text{ LI}$$

- Perform a 3rd-order polynomial regression and sketch the result of line fitting.
- Given a test point $\{x = 9\}$ predict y using the polynomial model.
- Compare this prediction with that of a linear regression.

(a) $X_{\text{train}} \Rightarrow P_{\text{train}} = \text{get_poly_data}(X_{\text{train}}, 3)$
 $\hat{w} = \text{LI}(P_{\text{train}}, y)$

(b) $X_{\text{test}} = (9) \rightarrow P_{\text{test}} = \begin{pmatrix} 1 & 9 & 9^2 & 9^3 \end{pmatrix}$ (2) $y = P_{\text{test}} \hat{w} = f(X_{\text{test}})$

(Polynomial Regression, 3D data, Python)

Question 3:

- Write down the expression for a 3rd order polynomial model having a 3-dimensional input.

$$y = w_{0,0,0} + w_{1,0,0}x_1 + w_{0,1,0}x_2 + w_{0,0,1}x_3 + w_{2,0,0}x_1^2 + w_{1,1,0}x_1x_2 + w_{1,0,1}x_1x_3 + w_{0,2,0}x_2^2 + \dots + w_{3,0,0}x_1^3 + \dots + w_{0,0,3}x_3^3$$

- Write down the P matrix for this polynomial given $X = \begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 1 \\ 1 & -1 & 1 \end{bmatrix}$. What to check?
- Given $y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$, can a unique solution be obtained in dual form? If so, proceed to solve it. $m = 2$, $d_p = 20$
- Given $y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$, can the primal ridge regression be applied to obtain a unique solution? If so, proceed to solve it. by right mcd $P^T P$ invertible / not invertible. But: $P^T P + \lambda I$ invertible.
-

(Binary Classification, Python)

Question 4:

Given the training data:

- $\{x = -1\} \rightarrow \{y = \text{class1}\}$
- $\{x = 0\} \rightarrow \{y = \text{class1}\}$
- $\{x = 0.5\} \rightarrow \{y = \text{class2}\}$
- $\{x = 0.3\} \rightarrow \{y = \text{class1}\}$
- $\{x = 0.8\} \rightarrow \{y = \text{class2}\}$

Predict the class label for $\{x = -0.1\}$ and $\{x = 0.4\}$ using linear regression with signum discrimination.

(1) $y = \begin{pmatrix} +1 \\ +1 \\ -1 \\ -1 \\ -1 \end{pmatrix}$

(2) Do lin. reg to get \hat{w} .
 $y\text{-pred} = X_{\text{test}} \hat{w} = \begin{pmatrix} 1 & -0.1 \\ 1 & 0.4 \end{pmatrix} \cdot \hat{w}$

(3) $\text{class-pred} = \text{sgn}(y\text{-pred})$

How many monomials?

to count monomials up to deg K from

$$\{x_1, \dots, x_d\}$$

e.g. to count

$$\{1, x_1, \dots, x_d, x_1^2, x_1 x_2, \dots, x_1 x_n, x_2^2, \dots, x_n^2\}$$

\Leftrightarrow counting

$$\{x_0^K, x_0^{K-1} x_1, \dots, x_0^{K-1} x_d, x_0^{K-2} x_1^2, \dots\}$$

and set $x_0 = 1$.

\Leftrightarrow choosing K items from a list of $d+1$ items with replacement $\{x_0 x_1, \dots, x_d\}$

Formula: $\binom{(d+1)+K-1}{K} = \binom{d+K}{K}$

$w_{i,j,k} \rightarrow$ weight attached to $x_1^i x_2^j x_3^k$

$$\binom{3+3}{3} = \binom{6}{3} = \frac{6 \times 5 \times 4}{3 \times 2} = 20$$

normally, $y \sim (m, 1)$ array
 now (one-hot-encoding), $Y \sim (m, c)$ array where c is no. of classes

(Multi-Category Classification, Python)

Question 5:

Given the training data:

$$\begin{aligned} \{x = -1\} &\rightarrow \{y = \text{class1}\} \\ \{x = 0\} &\rightarrow \{y = \text{class1}\} \\ \{x = 0.5\} &\rightarrow \{y = \text{class2}\} \\ \{x = 0.3\} &\rightarrow \{y = \text{class3}\} \\ \{x = 0.8\} &\rightarrow \{y = \text{class2}\} \end{aligned}$$

- (a) Predict the class label for $\{x = -0.1\}$ and $\{x = 0.4\}$ based on linear regression towards a one-hot encoded target.

$m = 5$ $d = 1 + 1 = 2$ $Y_{\text{pred}} = X_{\text{test-avg}} \cdot \hat{w} = \begin{pmatrix} 1 & -0.1 \\ 1 & 0.4 \end{pmatrix} \hat{w}$

- (b) Predict the class label for $\{x = -0.1\}$ and $\{x = 0.4\}$ using a polynomial model of 5th order and a one-hot encoded target.

$X_{\text{train}} \rightarrow P_{\text{train}}$
 \downarrow
 $\text{shape} =$
 \downarrow
 $f(x) =$
 $\left(\begin{matrix} 1 & x & x^2 & x^3 & x^4 & x^5 \end{matrix} \right) \begin{pmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \\ w_4 \\ w_5 \end{pmatrix}$
 $\rightarrow f(x) = P \cdot w$

(Multi-Category Classification, Python)

Question 6 (continued from Q3 of Tutorial 2):

Get the data set "from sklearn.datasets import load_iris". Use Python to perform the following tasks.

- Split the database into two sets: 74% of samples for training, and 26% of samples for testing. Hint: you might want to utilize from sklearn.model_selection import train_test_split for the splitting.
- Construct the target output using one-hot encoding.
- Perform a linear regression for classification (without inclusion of ridge, utilizing one-hot encoding for the learning target) and compute the number of test samples that are classified correctly.
- Using the same training and test sets as in above, perform a 2nd order polynomial regression for classification (again, without inclusion of ridge, utilizing one-hot encoding for the learning target) and compute the number of test samples that are classified correctly. Hint: you might want to use from sklearn.preprocessing import PolynomialFeatures for generation of the polynomial matrix.

Question 7

MCQ: there could be more than one answer. Given three samples of two-dimensional data points $X = \begin{bmatrix} 1 & 1 \\ 0 & 1 \\ 3 & 3 \end{bmatrix}$ with

corresponding target vector $y = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$. Suppose you want to use a full third-order polynomial model to fit these data.

Which of the following is/are true?

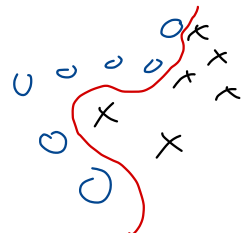
- The polynomial model has 10 parameters to learn T $d_p = 10$.
- The polynomial learning system is an under-determined one $T < d_p$.
- The learning of the polynomial model has infinite number of solutions T .
- The input matrix X has linearly dependent samples T .
- None of the above

Linearly dep.

$m > d$: "no soln".

$m < d$: under constrained.
 & inf. solutions.

no. of parameters = $\binom{d+k}{k}$
 $= \binom{3+2}{3} = \binom{5}{3}$
 $= \frac{5!}{3!(5-3)!} = \frac{5 \times 4}{2} = 10$



Question 8

MCQ: there could be more than one answer. Which of the following is/are true?

- a) The polynomial model can be used to solve problems with nonlinear decision boundary. **T.**
- b) The ridge regression cannot be applied to multi-target regression. **T.**
- c) The solution for learning feature \mathbf{X} with target \mathbf{y} based on linear ridge regression can be written as $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$ for $\lambda > 0$. As λ increases, $\hat{\mathbf{w}}^T \hat{\mathbf{w}}$ decreases. **T.**
- d) If there are four data samples with two input features each, the full second-order polynomial model is an over-determined system. **F.**

$$X = 4 \begin{pmatrix} 2 \\ 2 \end{pmatrix}$$

$$\left. \begin{matrix} m = 4 \\ d_p = 6 \end{matrix} \right\} \text{ () - determined -}$$

$$d = \binom{2+2}{2} = \frac{4 \cdot 3}{2} = 6.$$

→ "add λI ".
 → classically > 2 classes.
 ① OLS
 ② Linear/Adg regression.
 ③ PLS D.

$$\text{cost} = \text{MSE}(\text{pred}, \text{true}) + \lambda \|\mathbf{w}\|^2$$