

EE2211 Pre-Tutorial 6

Dr Feng LIN

feng_lin@nus.edu.sg



Agenda

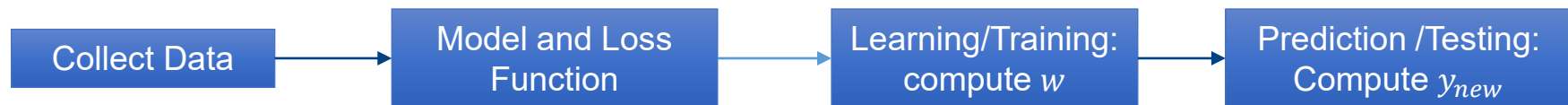
- Recap
- Self-learning
- Tutorial 6



Recap

- Linear Classification
 - Binary classification
 - Multi-category classification
- Ridge regression
 - Penalty term
 - Primal and dual forms
- Polynomial Regression
 - Nonlinear decision boundary

Linear Regression



$$\mathbf{X}\mathbf{w} = \mathbf{y}$$

$$\frac{1}{m} \sum_{i=1}^m (f_{\mathbf{w},b}(\mathbf{x}_i) - y_i)^2$$

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\hat{f}_{\mathbf{w}}(\mathbf{X}_{new}) = \mathbf{X}_{new} \hat{\mathbf{w}}$$

- \mathbf{X} : Samples
- \mathbf{y} : Target values

- Linear or Affine function
- Squared error loss function

- Check the invertibility
- Least square approximation (left-inverse)

- Prediction for new inputs
- Testing: Mean Squared Error (MSE)

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Linear Regression

For one sample: a linear model $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{x}^T \mathbf{w}$

scalar function

For m samples: $f_{\mathbf{w}}(\mathbf{X}) = \mathbf{X}\mathbf{w} = \mathbf{y}$

$$\mathbf{y} = \begin{bmatrix} \mathbf{x}_1^T \mathbf{w} \\ \vdots \\ \mathbf{x}_m^T \mathbf{w} \end{bmatrix} \quad \text{where} \quad \mathbf{x}_i^T = [1, x_{i,1}, \dots, x_{i,d}]$$
$$\mathbf{X} = \begin{bmatrix} 1 & x_{1,1} & \dots & x_{1,d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{m,1} & \dots & x_{m,d} \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} b \\ w_1 \\ \vdots \\ w_d \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}$$

Objective: $\sum_{i=1}^m (f_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2 = \mathbf{e}^T \mathbf{e} = (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$

Learning/training

when $\mathbf{X}^T \mathbf{X}$ is invertible

Least square solution: $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

Prediction/testing:

$\mathbf{y}_{new} = \hat{f}_{\mathbf{w}}(\mathbf{X}_{new}) = \mathbf{X}_{new} \hat{\mathbf{w}}$

Linear Regression

Learning (Training)

- Consider the set of feature vector \mathbf{x}_i and target output y_i indexed by $i = 1, \dots, m$, a linear model $f_w(\mathbf{x}) = \mathbf{x}^T \mathbf{w}$ can be stacked as

$$\begin{aligned}
 \underset{\text{Learning Model}}{f_w(\mathbf{X}) = \mathbf{X}\mathbf{w}} & \longleftrightarrow \underset{\text{Learning target vector}}{\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}} \\
 &= \begin{bmatrix} \mathbf{x}_1^T \mathbf{w} \\ \vdots \\ \mathbf{x}_m^T \mathbf{w} \end{bmatrix}
 \end{aligned}$$

where $\mathbf{x}_i^T \mathbf{w} = [1, x_{i,1}, \dots, x_{i,d}] \begin{bmatrix} b \\ w_1 \\ \vdots \\ w_d \end{bmatrix}$

Objective: $\sum_{i=1}^m (f_w(\mathbf{x}_i) - y_i)^2 = \mathbf{e}^T \mathbf{e} = (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$

$$\sum_{i=1}^m (f_w(\mathbf{x}_i) - y_i)^2 = \sum_{i=1}^m (x_i^T \mathbf{w} - y_i)^2 = (x_1^T \mathbf{w} - y_1)^2 + (x_2^T \mathbf{w} - y_2)^2 + \dots + (x_m^T \mathbf{w} - y_m)^2$$

$$= [x_1^T \mathbf{w} - y_1 \quad x_2^T \mathbf{w} - y_2 \quad \dots \quad x_m^T \mathbf{w} - y_m] \begin{bmatrix} x_1^T \mathbf{w} - y_1 \\ x_2^T \mathbf{w} - y_2 \\ \vdots \\ x_m^T \mathbf{w} - y_m \end{bmatrix}$$

$$= (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) = \mathbf{e}^T \mathbf{e}$$

Linear Classification

Linear Methods for Classification

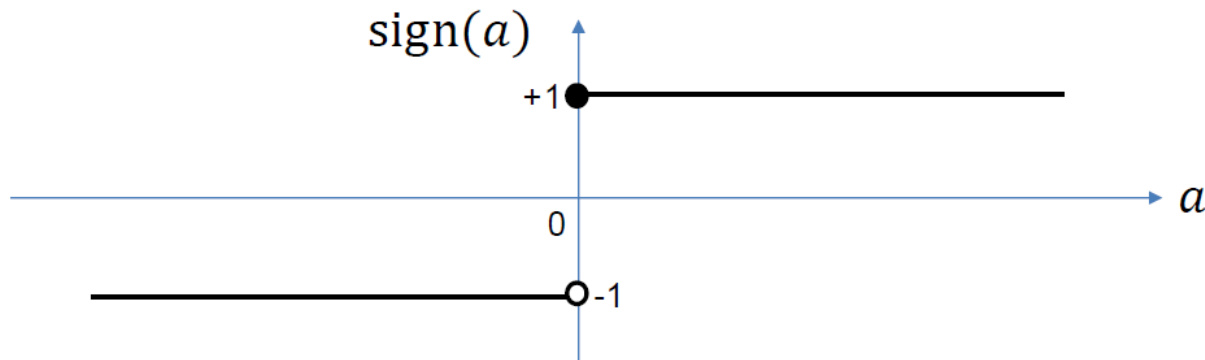
Binary Classification:

If $\mathbf{X}^T \mathbf{X}$ is invertible, then

Learning: $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$, $y_i \in \{-1, +1\}, i = 1, \dots, m$

Prediction: $\hat{f}_{\mathbf{w}}^c(\mathbf{x}_{new}) = \text{sign}(\mathbf{x}_{new}^T \hat{\mathbf{w}})$ for each row \mathbf{x}_{new}^T of \mathbf{X}_{new}

$\text{sign}(a) = +1$ for $a \geq 0$ and -1 for $a < 0$



Linear Classification

Linear Methods for Classification

Multi-Category Classification:

If $\mathbf{X}^T \mathbf{X}$ is invertible, then

Learning: $\hat{\mathbf{W}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$, $\mathbf{Y} \in \mathbb{R}^{m \times C}$

Prediction: $\hat{f}_{\mathbf{w}}^c(\mathbf{x}_{new}) = \arg \max_{k=1, \dots, C} \left(\mathbf{x}_{new}^T \hat{\mathbf{W}}(:, k) \right)$ for each \mathbf{x}_{new}^T of \mathbf{X}_{new}

Each row (of $i = 1, \dots, m$) in \mathbf{Y} has an **one-hot** encoding/assignment:

e.g., target for class-1 is labelled as $\mathbf{y}_i^T = [1, 0, 0, \dots, 0]$ for the i th sample,
target for class-2 is labelled as $\mathbf{y}_j^T = [0, 1, 0, \dots, 0]$ for the j th sample,
target for class-C is labelled as $\mathbf{y}_m^T = [0, 0, \dots, 0, 1]$ for the m th sample.

C

Ridge Regression

Recall Linear regression

Objective: $\hat{\mathbf{w}} = \operatorname{argmin} \sum_{i=1}^m (f_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2 = (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$

The learning computation: $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

We cannot guarantee that the matrix $\mathbf{X}^T \mathbf{X}$ is invertible

Ridge regression: shrinks the regression coefficients w by imposing a penalty on their size

Objective: $\hat{\mathbf{w}} = \operatorname{argmin} \sum_{i=1}^m (f_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2 + \lambda \sum_{j=1}^d w_j^2$
 $= \operatorname{argmin} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w}$

Here $\lambda \geq 0$ is a complexity parameter that controls the amount of shrinkage: the larger the value of λ , the greater the amount of shrinkage.

Note: m samples & d parameters

Ridge Regression

The learning computation: $\hat{\mathbf{w}} = (X^T X)^{-1} X^T \mathbf{y}$

$$(X^T X)^{-1} = \frac{1}{|X^T X|} (X^T X)^* \rightarrow \frac{1}{0} (X^T X)^* \Rightarrow \hat{\mathbf{w}} \rightarrow \infty$$

If $X^T X$ is not invertible, that means its determination is 0. This causes the denominator of $(X^T X)^{-1}$ to approach 0, which in turn **causes w to approach infinity**, making it impossible to fit the data well.

Ridge regression: shrinks the regression coefficients w by impose **penalty on their size**

$$\begin{aligned} \text{Objective: } \hat{\mathbf{w}} &= \operatorname{argmin} \sum_{i=1}^m (f_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2 + \lambda \sum_{j=1}^d w_j^2 \\ &= \operatorname{argmin} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w} \end{aligned}$$

$$\lambda \|\mathbf{w}\|^2$$

Regularization
or penalty term
or ridge term

Ridge Regression

Using a linear model:

$$\min_{\mathbf{w}} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w}$$

Solution:

$$\frac{\partial}{\partial \mathbf{w}} ((\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w}) = \mathbf{0}$$

$$\Rightarrow 2\mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{X}^T \mathbf{y} + 2\lambda \mathbf{w} = \mathbf{0}$$

$$\Rightarrow \mathbf{X}^T \mathbf{X} \mathbf{w} + \lambda \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

$$\Rightarrow (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

where \mathbf{I} is the $d \times d$ identity matrix

Here on, we shall focus on single column of output \mathbf{y} in derivations in the sequel

Learning: $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$

Page 10, Lecture 5

$$\frac{d\mathbf{A}\mathbf{x}}{d\mathbf{x}} = \mathbf{A}$$

$$\frac{d(\mathbf{b}^T \mathbf{x})}{d\mathbf{x}} = \mathbf{b} \quad \frac{d(\mathbf{y}^T \mathbf{A}\mathbf{x})}{d\mathbf{x}} = \mathbf{A}^T \mathbf{y}$$

$$\frac{d(\mathbf{x}^T \mathbf{A}\mathbf{x})}{d\mathbf{x}} = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}$$

Ridge Regression

Ridge Regression in Primal Form (when $m > d$)

$(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})$ is invertible for $\lambda > 0$,

Learning: $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$

Prediction: $\hat{f}_{\mathbf{w}}(\mathbf{X}_{new}) = \mathbf{X}_{new} \hat{\mathbf{w}}$

Ridge Regression in Dual Form (when $m < d$)

$(\mathbf{X} \mathbf{X}^T + \lambda \mathbf{I})$ is invertible for $\lambda > 0$,

Learning: $\hat{\mathbf{w}} = \mathbf{X}^T (\mathbf{X} \mathbf{X}^T + \lambda \mathbf{I})^{-1} \mathbf{y}$

Prediction: $\hat{f}_{\mathbf{w}}(\mathbf{X}_{new}) = \mathbf{X}_{new} \hat{\mathbf{w}}$

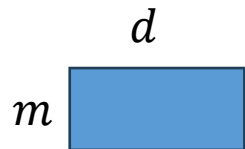
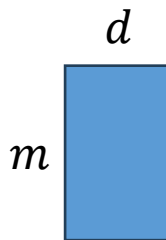
Linear and Ridge Regression

	Linear Regression	Ridge Regression
Over-determined system ($m > d$)	Left inverse $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$	Primal Form $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$
Under-determined system ($m < d$)	Right inverse $\hat{\mathbf{w}} = \mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{y}$	Dual Form $\hat{\mathbf{w}} = \mathbf{X}^T (\mathbf{X} \mathbf{X}^T + \lambda \mathbf{I})^{-1} \mathbf{y}$

Noted: 1) The primal form can be used to solve under-determined system, but it is better suited for over-determined system. 2) The dual form of ridge regression is often more computationally efficient in under-determined system than the primal form.

When to Use Primal vs. Dual Form?

Ridge Regression	Matrix Inversion Size	Best for
<p>Primal Form</p> $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$	$d \times d$	<p>Over-determined system ($m > d$)</p>
<p>Dual Form</p> $\hat{\mathbf{w}} = \mathbf{X}^T (\mathbf{X} \mathbf{X}^T + \lambda \mathbf{I})^{-1} \mathbf{y}$	$m \times m$	<p>Under-determined system ($m < d$)</p>



Polynomial Regression

Motivation: nonlinear decision surface

- Based on the sum of products of the variables
- E.g. when the input dimension is $d=2$,

a polynomial function of degree = 2 is:

$$f_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1 x_1 + w_2 x_2 + w_{12} x_1 x_2 + w_{11} x_1^2 + w_{22} x_2^2.$$

XOR problem

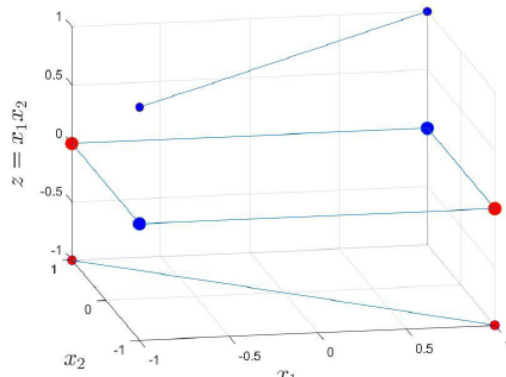
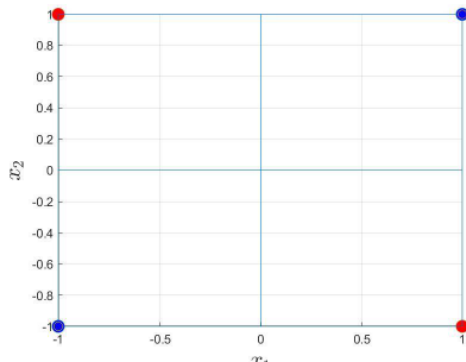
$$\mathbf{x}_1 = [+1 \ +1]^\top \quad y_1 = +1$$

$$\mathbf{x}_2 = [-1 \ +1]^\top \quad y_2 = -1$$

$$\mathbf{x}_3 = [+1 \ -1]^\top \quad y_3 = -1$$

$$\mathbf{x}_4 = [-1 \ -1]^\top \quad y_4 = +1$$

$$f_{\mathbf{w}}(\mathbf{x}) = x_1 x_2$$



Polynomial Regression

Motivation: Nonlinear Prediction

E.g. predicting the price of the house. Suppose you have two features:

- x_1 : the frontage of house (the width of the property)
- x_2 : the depth of the house.

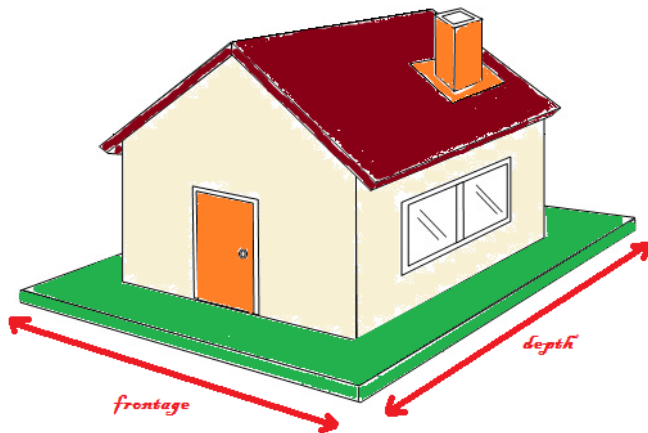
We might build a linear regression model like this

$$f_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2$$

If we want to predict house prices, we might focus on the house or land **area** as key factors and create a new feature accordingly.

$$f_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2 + w_{12}x_1x_2$$

Aera



Polynomial Regression

Polynomial Expansion

- The linear model $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{x}^T \mathbf{w}$ can be written as

$$\begin{aligned} f_{\mathbf{w}}(\mathbf{x}) &= \mathbf{x}^T \mathbf{w} \\ &= \sum_{i=0}^d x_i w_i, \quad x_0 = 1 \\ &= w_0 + \sum_{i=1}^d x_i w_i. \end{aligned}$$

- By including additional terms involving the products of pairs of components of \mathbf{x} , we obtain a quadratic model:

$$f_{\mathbf{w}}(\mathbf{x}) = w_0 + \sum_{i=1}^d w_i x_i + \sum_{i=1}^d \sum_{j=1}^d w_{ij} x_i x_j.$$

2nd order: $f_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1 x_1 + w_2 x_2 + w_{12} x_1 x_2 + w_{11} x_1^2 + w_{22} x_2^2$

3rd order: $f_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1 x_1 + w_2 x_2 + w_{12} x_1 x_2 + w_{11} x_1^2 + w_{22} x_2^2 + \sum_{i=1}^d \sum_{j=1}^d \sum_{k=1}^d w_{ijk} x_i x_j x_k, \quad d = 2$

n: number of variables
r: number of order

$$C(n, r) = \frac{(n + r - 1)!}{r! (n - 1)!}$$

$$C(3, 0) = \frac{(3 + 0 - 1)!}{0! (3 - 1)!} = 1$$

$$C(3, 1) = \frac{(3 + 1 - 1)!}{1! (3 - 1)!} = 3$$

$$C(3, 2) = \frac{(3 + 2 - 1)!}{2! (3 - 1)!} = 6$$

$$C(3, 3) = \frac{(3 + 3 - 1)!}{3! (3 - 1)!} = 10$$

Polynomial Regression

Generalized Linear Discriminant Function

$$f_{\mathbf{w}}(\mathbf{x}) = w_0 + \sum_{i=1}^d w_i x_i + \sum_{i=1}^d \sum_{j=1}^d w_{ij} x_i x_j + \sum_{i=1}^d \sum_{j=1}^d \sum_{k=1}^d w_{ijk} x_i x_j x_k + \dots$$

$$f_{\mathbf{w}}(\mathbf{x}) = \mathbf{P}\mathbf{w} \quad (\text{Note: } \mathbf{P} \triangleq \mathbf{P}(\mathbf{X}) \text{ for symbol simplicity})$$

$$= \begin{bmatrix} \mathbf{p}_1^T \mathbf{w} \\ \vdots \\ \mathbf{p}_m^T \mathbf{w} \end{bmatrix}$$

$$\text{where } \mathbf{p}_l^T \mathbf{w} = [1, x_{l,1}, \dots, x_{l,d}, \dots, x_{l,i}x_{l,j}, \dots, x_{l,i}x_{l,j}x_{l,k}, \dots]$$

$l = 1, \dots, m$; d denotes the dimension of input features; m denotes the number of samples

$$\begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \\ \vdots \\ w_{ij} \\ \vdots \\ w_{ijk} \\ \vdots \end{bmatrix}$$

Polynomial Regression

Ridge Regression in Primal Form ($m > d$)

For $\lambda > 0$,

Learning: $\hat{\mathbf{w}} = (\mathbf{P}^T \mathbf{P} + \lambda \mathbf{I})^{-1} \mathbf{P}^T \mathbf{y}$

Prediction: $\hat{f}_{\mathbf{w}}(\mathbf{P}(\mathbf{X}_{new})) = \mathbf{P}_{new} \hat{\mathbf{w}}$

Ridge Regression in Dual Form ($m < d$)

For $\lambda > 0$,

Learning: $\hat{\mathbf{w}} = \mathbf{P}^T (\mathbf{P} \mathbf{P}^T + \lambda \mathbf{I})^{-1} \mathbf{y}$

Prediction: $\hat{f}_{\mathbf{w}}(\mathbf{P}(\mathbf{X}_{new})) = \mathbf{P}_{new} \hat{\mathbf{w}}$

Note: Change \mathbf{X} to \mathbf{P} with reference to slides 15/16; m & d refers to the size of \mathbf{P} (not \mathbf{X})



THANK YOU