

INTRO TO DATA SCIENCE

LECTURE 6: REGRESSION & REGULARIZATION

LAST TIME:

I. WHAT IS MACHINE LEARNING?

II. MACHINE LEARNING PROBLEMS

III. CLASSIFICATION PROBLEMS

IV. KNN CLASSIFICATION (& EXERCISES)

QUESTIONS?

5. INTRO TO MACHINE LEARNING & KNN (LAST WEEK)

6. LINEAR REGRESSION & LINEAR ALGEBRA (TODAY)

7. REGRESSION & REGULARIZATION

8. STATISTICS & BAYES

9. DECISION TREES

10. RECAP SUPERVISED LEARNING

0. PRESENTATIONS ASSIGNMENT DATA EXPLORATION

I. LINEAR ALGEBRA & NUMPY

II. INTRO TO REGRESSION

III. MATH BEHIND THE SCENES (OPTIONAL)

INTRO TO DATA SCIENCE

I. LINEAR ALGEBRA

II. LINEAR REGRESSION

	<i>continuous</i>	<i>categorical</i>
<i>supervised</i>	<i>regression</i>	<i>classification</i>
<i>unsupervised</i>	<i>dimension reduction</i>	<i>clustering</i>

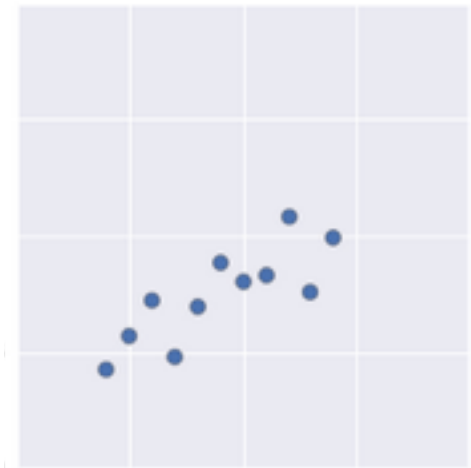
Q: What is a regression model?

*Q: What is a **regression** model?*

A: A functional relationship between input & response variables

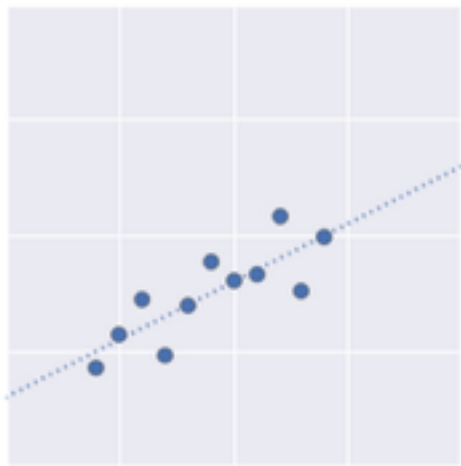
*Q: What is a **regression model**?*

A: A functional relationship between input & response variables



*Q: What is a **regression model**?*

A: A functional relationship between input & response variables



*Q: What is a **regression model**?*

A: A functional relationship between input & response variables

*The **simple linear regression model** captures a linear relationship between a single input variable x and a response variable y :*

*Q: What is a **regression model**?*

A: A functional relationship between input & response variables

*The **simple linear regression model** captures a linear relationship between a single input variable x and a response variable y :*

$$y = \alpha + \beta x + \varepsilon$$

Q: What do the terms in this model mean?

$$y = \alpha + \beta x + \varepsilon$$

Q: What do the terms in this model mean?

$$y = \alpha + \beta x + \varepsilon$$

y = response variable *(the one we want to predict)*

Q: What do the terms in this model mean?

$$y = \alpha + \beta x + \varepsilon$$

y = response variable *(the one we want to predict)*

x = input variable *(the one we use to train the model)*

Q: What do the terms in this model mean?

$$y = \alpha + \beta x + \varepsilon$$

y = response variable (*the one we want to predict*)

x = input variable (*the one we use to train the model*)

α = intercept (*where the line crosses the y -axis*)

Q: What do the terms in this model mean?

$$y = \alpha + \beta x + \varepsilon$$

y = response variable *(the one we want to predict)*

x = input variable *(the one we use to train the model)*

α = intercept *(where the line crosses the y -axis)*

β = regression coefficient *(the model “parameter”)*

Q: What do the terms in this model mean?

$$y = \alpha + \beta x + \varepsilon$$

y = **response variable** (*the one we want to predict*)

x = **input variable** (*the one we use to train the model*)

α = **intercept** (*where the line crosses the y -axis*)

β = **regression coefficient** (*the model “parameter”*)

ε = **residual** (*the prediction error*)

$$y = \alpha + \beta x + \varepsilon$$

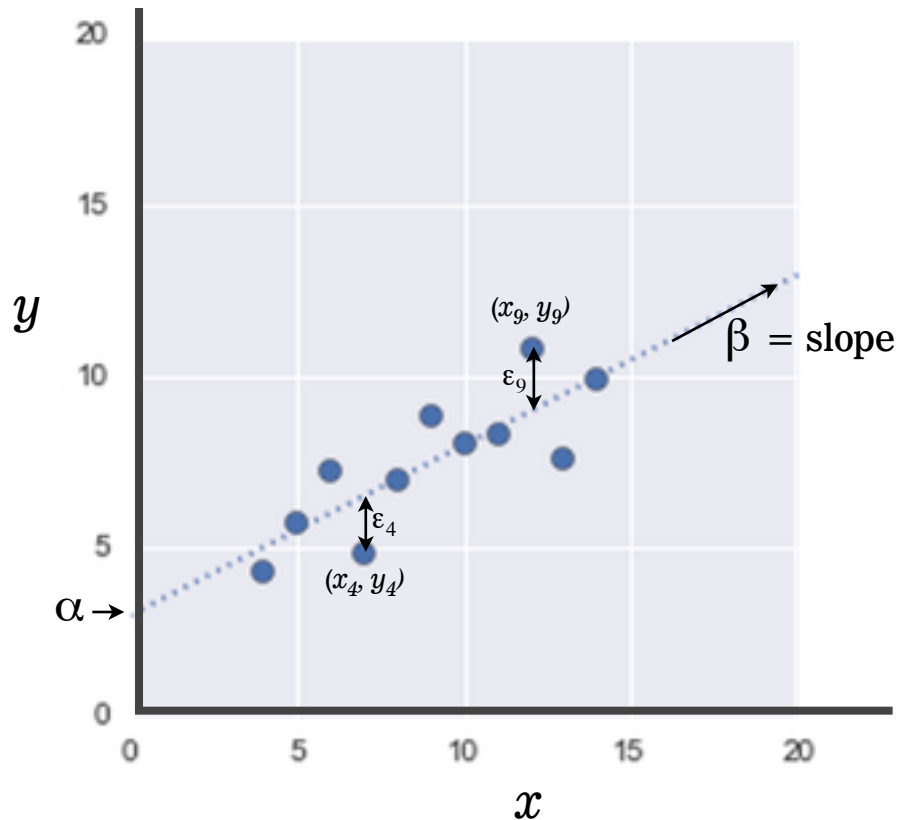
y = response variable

x = input variable

α = intercept

β = regression coefficient

ε = residual (*the error*)



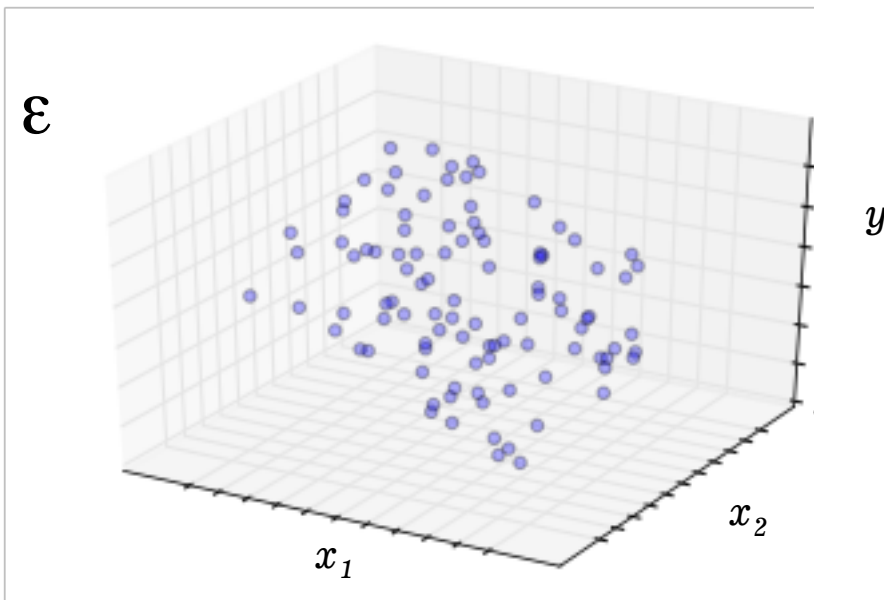
We can extend this model to several input variables, giving us the
multiple linear regression model:

We can extend this model to several input variables, giving us the multiple linear regression model:

$$y = \alpha + \beta_1 x_1 + \dots + \beta_n x_n + \varepsilon$$

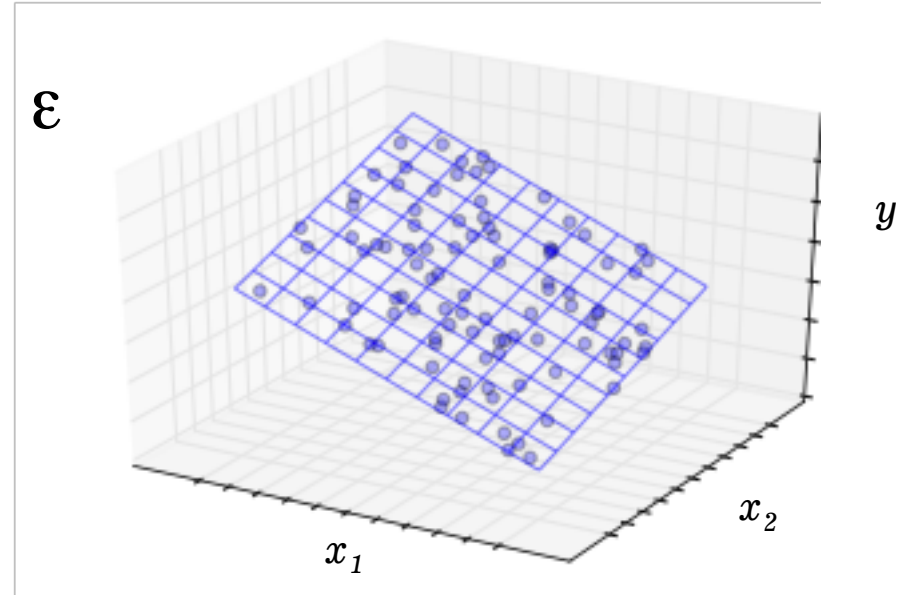
We can extend this model to several input variables, giving us the multiple linear regression model:

$$y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$$



We can extend this model to several input variables, giving us the multiple linear regression model:

$$y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$$



Q: How do we fit a regression model to a dataset?

Q: How do we fit a regression model to a dataset?

A: In theory, minimize the sum of the squared residuals (OLS).

Q: How do we fit a regression model to a dataset?

A: In theory, minimize the sum of the squared residuals (OLS).

In practice, any respectable piece of software will do this for you.

Q: How do we fit a regression model to a dataset?

A: In theory, minimize the sum of the squared residuals (OLS).

In practice, any respectable piece of software will do this for you.

But again, if you get serious about regression, you should learn how this works!

QUESTION

***HOW
DO YOU
MEASURE
THE
QUALITY?***

<i>supervised</i> <i>unsupervised</i>	<i>making predictions</i> <i>extracting structure</i>
--	--

supervised

test out your predictions

	<i>continuous</i>	<i>categorical</i>
<i>supervised</i>	<i>regression</i> <div>R^2 mean absolute error mean squared error</div>	<i>classification</i> <div>Accuracy (% correct predictions) and other metrics</div>

III: MATH BEHIND REGRESSION

Q: How do we fit a regression model to a dataset?

Given the multiple linear regression model:

$$y = \alpha + \beta_1 x_1 + \dots + \beta_n x_n + \varepsilon$$

Given the multiple linear regression model:

$$y = \alpha + \beta_1 x_1 + \dots + \beta_n x_n + \varepsilon$$

$$\hat{y} = \alpha + \beta_1 x_1 + \dots + \beta_n x_n$$

The residual (error) is equal to the observed y minus the predicted \hat{y}

$$\begin{aligned}\varepsilon &= -y + \alpha + \beta_1 x_1 + \dots + \beta_n x_n \\ &= \hat{y} - y\end{aligned}$$

*Define a **cost function** J of the parameters α and β s*

$$J(\alpha, \beta) = \sum (-y + \alpha + \beta_1 x_1 + \dots + \beta_n x_n)^2$$

which sums the squares of all prediction errors

Then we're looking for those β s where J has its minimum:

$$\min_{\alpha, \beta} J = \min_{\alpha, \beta} \sum (-y + \alpha + \beta_1 x_1 + \dots + \beta_n x_n)^2$$

*Define a **cost function** J of the parameters β*

$$\min_{\alpha, \beta} J = \min_{\alpha, \beta} \sum (-y + \alpha + \beta_1 x_1 + \dots + \beta_n x_n)^2$$

and find where J has its minimum

This is called the Ordinary Least Squares (OLS) method

Let's simplify notation using linear algebra

Given the multiple linear regression model:

$$y = \alpha + \beta_1 x_1 + \dots + \beta_n x_n + \varepsilon$$

Given the multiple linear regression model:

$$\begin{array}{l} N \text{ samples} \left\{ \begin{array}{l} y_1 = \alpha + \beta_1 x_{11} + \dots + \beta_n x_{1n} + \varepsilon_1 \\ y_2 = \alpha + \beta_1 x_{21} + \dots + \beta_n x_{2n} + \varepsilon_2 \\ y_3 = \alpha + \beta_1 x_{31} + \dots + \beta_n x_{3n} + \varepsilon_3 \\ y_4 = \alpha + \beta_1 x_{41} + \dots + \beta_n x_{4n} + \varepsilon_4 \end{array} \right. \end{array}$$

Given the multiple linear regression model:

The diagram illustrates the multiple linear regression model for N samples and n features. It consists of four equations, one for each sample, grouped by a large left curly brace labeled "N samples". Each equation is of the form $y_i = \alpha + \beta_1 x_{i1} + \dots + \beta_n x_{in} + \epsilon_i$. A small bottom curly brace under the y_i terms is labeled "labels". A large bottom curly brace under the feature terms $\beta_1 x_{i1} + \dots + \beta_n x_{in}$ is labeled "n features".

$$\begin{array}{l} \left. \begin{array}{l} y_1 = \alpha + \beta_1 x_{11} + \dots + \beta_n x_{1n} + \epsilon_1 \\ y_2 = \alpha + \beta_1 x_{21} + \dots + \beta_n x_{2n} + \epsilon_2 \\ y_3 = \alpha + \beta_1 x_{31} + \dots + \beta_n x_{3n} + \epsilon_3 \\ y_4 = \alpha + \beta_1 x_{41} + \dots + \beta_n x_{4n} + \epsilon_4 \end{array} \right\} \begin{array}{l} N \text{ samples} \\ \text{labels} \end{array} \end{array}$$

$n \text{ features}$

Given the multiple linear regression model:

The diagram illustrates the multiple linear regression model for N samples. It shows four equations, each representing a sample. The first equation is $y_1 = \alpha + \beta_1 x_{11} + \dots + \beta_n x_{1n} + \epsilon_1$. The second is $y_2 = \alpha + \beta_1 x_{21} + \dots + \beta_n x_{2n} + \epsilon_2$. The third is $y_3 = \alpha + \beta_1 x_{21} + \dots + \beta_n x_{3n} + \epsilon_3$. The fourth is $y_4 = \alpha + \beta_1 x_{31} + \dots + \beta_n x_{4n} + \epsilon_4$. A large left curly brace groups these equations and is labeled "N samples". A bottom curly brace under the feature terms $\beta_1 x_{11} + \dots + \beta_n x_{1n}$ in the first equation is labeled "n features". A bottom curly brace under the error terms $\epsilon_1, \epsilon_2, \epsilon_3, \epsilon_4$ is labeled "residuals (errors)". A yellow oval highlights the error terms.

$$\begin{aligned} y_1 &= \alpha + \beta_1 x_{11} + \dots + \beta_n x_{1n} + \epsilon_1 \\ y_2 &= \alpha + \beta_1 x_{21} + \dots + \beta_n x_{2n} + \epsilon_2 \\ y_3 &= \alpha + \beta_1 x_{21} + \dots + \beta_n x_{3n} + \epsilon_3 \\ y_4 &= \alpha + \beta_1 x_{31} + \dots + \beta_n x_{4n} + \epsilon_4 \end{aligned}$$

N samples

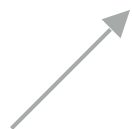
labels

n features

residuals (errors)

We can summarize all samples in vectors

$$\mathbf{y} = \alpha + \beta_1 \mathbf{x}_1 + \dots + \beta_n \mathbf{x}_n + \boldsymbol{\varepsilon}$$



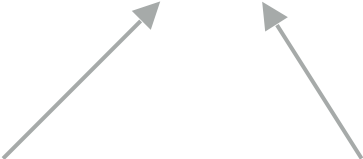
*\mathbf{y} , \mathbf{x} and $\boldsymbol{\varepsilon}$ are vectors of dimension N ,
where N is the number of samples (observations)*

And we can summarize all features ...

$$\mathbf{y} = \alpha + \beta_1 \mathbf{x}_1 + \dots + \beta_n \mathbf{x}_n + \boldsymbol{\varepsilon}$$



And we can summarize all features in a matrix \mathbf{X}

$$\mathbf{y} = \alpha + \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$


*\mathbf{X} is an $N \times n$ -matrix,
with N rows for each observation (samples)
and n columns for each feature*

*$\boldsymbol{\beta}$ is a n -dimensional vector,
with coefficients for each feature*

So we can write the linear regression, using linear algebra, as

$$\mathbf{y} = \alpha + \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

So we can write the linear regression, using linear algebra, as

$$\mathbf{y} = \alpha + \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

y = response variable <i>(the one we want to predict)</i>	<i>N-dim vector</i>
X = input variable <i>(the one we use to train the model)</i>	<i>N×n-matrix</i>
α = intercept <i>(where the line crosses the y-axis)</i>	<i>scalar</i>
β = regression coefficient <i>(the model “parameter”)</i>	<i>n-dim vector</i>
ε = residual <i>(the prediction error)</i>	<i>N-dim vector</i>

So we can write the linear regression, using linear algebra, as

$$\mathbf{y} = \alpha + \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

And hence the cost function as the norm of the residual vector $\boldsymbol{\varepsilon}$

$$J(\alpha, \boldsymbol{\beta}) = |-\mathbf{y} + \alpha + \mathbf{X}\boldsymbol{\beta}|$$

So we can write the linear regression, using linear algebra, as

$$\mathbf{y} = \alpha + \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

And hence the cost function as the norm of the residual vector $\boldsymbol{\varepsilon}$

$$\begin{aligned} J(\alpha, \boldsymbol{\beta}) &= |-\mathbf{y} + \alpha + \mathbf{X}\boldsymbol{\beta}| \\ &= \sum_{i=1}^N (-y_i + \alpha + \mathbf{x}_i\boldsymbol{\beta})^2 \end{aligned}$$

So we can write the linear regression, using linear algebra, as

$$\mathbf{y} = \alpha + \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

And hence the cost function as the norm of the residual vector $\boldsymbol{\varepsilon}$

$$J(\alpha, \boldsymbol{\beta}) = |-\mathbf{y} + \alpha + \mathbf{X}\boldsymbol{\beta}|$$

which is a function $\mathbb{R}^{n+1} \rightarrow \mathbb{R}$, for which we want to know the minimum

Often an artificial 0th column of 1s is added to \mathbf{X} , such that we get

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

and $\alpha = \beta_0$, which is often easier in notation and computations

$$J(\boldsymbol{\beta}) = |-\mathbf{y} + \mathbf{X}\boldsymbol{\beta}|$$

*\mathbf{X} is now an $N \times (n+1)$ -matrix
 $\mathbf{X} = (1, \mathbf{x}_1, \mathbf{x}_2, \dots)$*

which is a function $\mathbb{R}^{n+1} \rightarrow \mathbb{R}$, for which we want to know the minimum

The OLS method has a closed-form solution for the parameter β

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

The OLS method has a closed-form solution for the parameter β

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

*Generally, though, you do **not** have a closed-form solution to find the minimum of the cost function. In that case, the **gradient descent** algorithm will help you out. (To Be Continued.)*

QUESTION

***HOW
DO YOU
MEASURE
THE
QUALITY?***

*An intuitive measure is the **mean absolute error***

$$\text{MAE} = \frac{1}{N} \sum |\varepsilon| = \frac{1}{N} \sum |y - \hat{y}|$$

*An intuitive measure is the **mean absolute error***

$$\text{MAE} = \frac{1}{N} \sum |\epsilon| = \frac{1}{N} \sum |y - \hat{y}|$$

*Another widely used one is the **mean squared error***

$$\text{MSE} = \frac{1}{N} \sum \epsilon^2 = \frac{1}{N} \sum (y - \hat{y})^2$$

*An intuitive measure is the **mean absolute error***

$$\text{MAE} = \frac{1}{N} \sum |\varepsilon| = \frac{1}{N} \sum |y - \hat{y}|$$

*Another widely used one is the **mean squared error***

$$\text{MSE} = \frac{1}{N} \sum \varepsilon^2 = \frac{1}{N} \sum (y - \hat{y})^2$$

Note that the OLS method minimizes the MSE

The most commonly used measure is R^2 , which is a ratio of variance

$$\text{MSE} = \frac{1}{N} \sum (y - \hat{y})^2 \quad \text{mean squared error}$$

The most commonly used measure is R^2 , which is a ratio of variance

$$\text{MSE} = \frac{1}{N} \sum (y - \hat{y})^2 \quad \text{mean squared error}$$

$$\text{Var } y = \frac{1}{N} \sum (y - \bar{y})^2 \quad \text{variance of observed data}$$

The most commonly used measure is R^2 , which is a ratio of variance

$$R^2 = 1 - \frac{\frac{1}{N} \sum (y - \hat{y})^2}{\frac{1}{N} \sum (y - \bar{y})^2}$$

mean squared error

variance of observed data

The most commonly used measure is R^2 , which is a ratio of variance

$$R^2 = 1 - \frac{\sum (y - \hat{y})^2}{\sum (y - \bar{y})^2}$$

The most commonly used measure is R^2 , which is a ratio of variance

$$R^2 = 1 - \frac{\sum (y - \hat{y})^2}{\sum (y - \bar{y})^2}$$

$R^2 = 1$ iff all errors are zero

The most commonly used measure is R^2 , which is a ratio of variance

$$R^2 = 1 - \frac{\sum (y - \hat{y})^2}{\sum (y - \bar{y})^2}$$

$R^2 = 1$ iff all errors are zero

$R^2 < 0$ is possible (very bad)

INTRO TO DATA SCIENCE

DISCUSSION