Machine Learning Regression Methods

# FAST DISCOVERING THE FUTURE

- ✓ Ingredients of Machine Learning
- ✓ Classification Basics
- ✓ Basic Linear Classifier
- ✓ K-Nearest Neighbours Classifier
- ✓ Naive Bayes Classifier
- ✓ Linear and Quadratic Discriminant Analysis
- ✓ Support Vector Machines (SVM)
- ✓ Decision Trees
- ✓ Ensemble Methods (Bagging, Weighted Voting, Stacking)

- Main concepts in regression
- Linear Regression
- Ordinary Least Squares (OLS)

• Regression task is the same as classification task, except that we must predict a continuous variable (instead of a categorical class label)

- Regression task is the same as classification task, except that we must predict a continuous variable (instead of a categorical class label)
- For example:

- Regression task is the same as classification task, except that we must predict a continuous variable (instead of a categorical class label)
- For example:
  - predict the salary given the info about a person

- Regression task is the same as classification task, except that we must predict a continuous variable (instead of a categorical class label)
- For example:
  - predict the salary given the info about a person
  - predict the risk of a policyholder for insurance

- Regression task is the same as classification task, except that we must predict a continuous variable (instead of a categorical class label)
- For example:
  - predict the salary given the info about a person
  - predict the risk of a policyholder for insurance
  - predict the expected number of days that a patient will stay in a hospital
  - ...



- X input space (set of all possible instances)
- $\mathbb{Y}$  output space (all possible labels)
- $f: \mathbb{X} \to \mathbb{Y}$  any such function is a classifier
- $\mathbf{x} \in \mathbb{X}$  instance
- $y \in \mathbb{Y}$  actual / true label of instance  $\mathbf{x}$
- $\hat{y} = f(\mathbf{x})$  predicted label of instance  $\mathbf{x}$

- X input space (set of all possible instances)
- Y output space (all possible labels real numbers)
- $f: \mathbb{X} \to \mathbb{Y}$  any such function is a classifier regression model
- $\mathbf{x} \in \mathbb{X}$  instance
- $y \in \mathbb{Y}$  actual / true <del>label</del> target value of instance  $\mathbf{x}$
- $\hat{y} = f(\mathbf{x})\text{-}$  predicted label target value of instance  $\mathbf{x}$

• Suppose that there exists an actual / true function, mapping the features to target variable  $f^*: X \to \mathbb{R}$ 

- Suppose that there exists an actual / true function, mapping the features to target variable  $f^*:\mathbb{X}\to\mathbb{R}$
- In regression the task is to learn a function approximator  $\hat{f}:\mathbb{X}\to\mathbb{R}$  such that  $\hat{f}\approx f^*$

- Suppose that there exists an actual / true function, mapping the features to target variable  $f^*:\mathbb{X}\to\mathbb{R}$
- In regression the task is to learn a function approximator  $\hat{f}:\mathbb{X}\to\mathbb{R}$  such that  $\hat{f}\approx f^*$
- For this we are given training data  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n) \in \mathbb{X} \times \mathbb{R}$

## • Do we want to learn $\hat{f}$ such that: $\hat{f}(\mathbf{x}_1) \approx y_1, \dots, \hat{f}(\mathbf{x}_n) \approx y_n$ ?

- Do we want to learn  $\hat{f}$  such that:  $\hat{f}(\mathbf{x}_1) \approx y_1, \dots, \hat{f}(\mathbf{x}_n) \approx y_n$  ?
- No! This would mean good predictions on training data, which is not our main goal!

- Do we want to learn  $\hat{f}$  such that:  $\hat{f}(\mathbf{x}_1) \approx y_1, \dots, \hat{f}(\mathbf{x}_n) \approx y_n$  ?
- No! This would mean good predictions on training data, which is not our main goal!
- We want to predict well on (future) test data!

- Do we want to learn  $\hat{f}$  such that:  $\hat{f}(\mathbf{x}_1) \approx y_1, \dots, \hat{f}(\mathbf{x}_n) \approx y_n$  ?
- No! This would mean good predictions on training data, which is not our main goal!
- We want to predict well on (future) test data!
- On any future instance  $\mathbf{X}\in\mathbb{X}$  with true (hidden) target Y we want  $\widehat{f}(\mathbf{X})\approx Y$

- Do we want to learn  $\hat{f}$  such that:  $\hat{f}(\mathbf{x}_1) \approx y_1, \dots, \hat{f}(\mathbf{x}_n) \approx y_n$  ?
- No! This would mean good predictions on training data, which is not our main goal!
- We want to predict well on (future) test data!
- On any future instance  $\mathbf{X}\in\mathbb{X}$  with true (hidden) target Y we want  $\widehat{f}(\mathbf{X})\approx Y$
- What does this really mean?

## Learning Problem in Regression



	lation
I A.) I	ылоп

#### Definition

For X and Y continuous random variables, the conditional expectation is

$$\mathbb{E}(X|Y) = \int_{x \in \mathcal{X}} x p(x|y) dx,$$

where  $p(x|y) = \frac{p(x,y)}{p(y)}$  is the conditional density function of X given Y.

#### Definition

For X and Y continuous random variables, the conditional expectation is

$$\mathbb{E}(X|Y) = \int_{x \in \mathcal{X}} x p(x|y) dx,$$

where  $p(x|y) = \frac{p(x,y)}{p(y)}$  is the conditional density function of X given Y.

• When choosing a particular estimate function  $f({\bf x})$  for the true value y we encounter some loss  $L(y,f({\bf x}))$ 

#### Definition

For X and Y continuous random variables, the conditional expectation is

$$\mathbb{E}(X|Y) = \int_{x \in \mathcal{X}} x p(x|y) dx,$$

where  $p(x|y) = \frac{p(x,y)}{p(y)}$  is the conditional density function of X given Y.

- When choosing a particular estimate function  $f({\bf x})$  for the true value y we encounter some loss  $L(y,f({\bf x}))$
- The average or expected loss is given by

$$\mathbb{E}(L) = \iint L(y, f(\mathbf{x})) p(\mathbf{x}, y) d\mathbf{x} dy$$

• The average or expected loss is given by  $\mathbb{E}(L)=\int\int L(y,f(\mathbf{x}))p(\mathbf{x},y)d\mathbf{x}dy$ 

- The average or expected loss is given by  $\mathbb{E}(L) = \iint L(y, f(\mathbf{x})) p(\mathbf{x}, y) d\mathbf{x} dy$
- A common choice of loss function in regression is the squared loss  $L(y, f(\mathbf{x})) = (f(\mathbf{x}) y)^2$ . The average loss will be

$$\mathbb{E}(L) = \iint (f(\mathbf{x}) - y)^2 p(\mathbf{x}, y) d\mathbf{x} dy$$

- The average or expected loss is given by  $\mathbb{E}(L) = \iint L(y, f(\mathbf{x})) p(\mathbf{x}, y) d\mathbf{x} dy$
- A common choice of loss function in regression is the squared loss  $L(y, f(\mathbf{x})) = (f(\mathbf{x}) y)^2$ . The average loss will be

$$\mathbb{E}(L) = \iint (f(\mathbf{x}) - y)^2 p(\mathbf{x}, y) d\mathbf{x} dy$$

• Our goal is to choose  $f(\mathbf{x})$  such that  $\mathbb{E}(L)$  will be minimized:

$$\frac{\partial \mathbb{E}(L)}{\partial f(\mathbf{x})} = 2 \int (f(\mathbf{x}) - y) p(\mathbf{x}, y) dy = 0$$

- The average or expected loss is given by  $\mathbb{E}(L) = \iint L(y, f(\mathbf{x})) p(\mathbf{x}, y) d\mathbf{x} dy$
- A common choice of loss function in regression is the squared loss  $L(y, f(\mathbf{x})) = (f(\mathbf{x}) y)^2$ . The average loss will be

$$\mathbb{E}(L) = \iint (f(\mathbf{x}) - y)^2 p(\mathbf{x}, y) d\mathbf{x} dy$$

• Our goal is to choose  $f(\mathbf{x})$  such that  $\mathbb{E}(L)$  will be minimized:

$$\frac{\partial \mathbb{E}(L)}{\partial f(\mathbf{x})} = 2 \int (f(\mathbf{x}) - y)p(\mathbf{x}, y)dy = 0$$

• Solving for  $f(\mathbf{x})$  gives

$$f(\mathbf{x}) = \frac{\int y p(\mathbf{x}, y) dy}{p(\mathbf{x})} = \int y p(y|\mathbf{x}) dy = \mathbb{E}(Y|X)$$

$$f^*(\mathbf{X}) = \mathbb{E}(Y|\mathbf{X})$$

$$f^*(\mathbf{X}) = \mathbb{E}(Y|\mathbf{X})$$

• Hence, for any  ${\bf X}$  we want to estimate the expected value of the corresponding Y

$$f^*(\mathbf{X}) = \mathbb{E}(Y|\mathbf{X})$$

- Hence, for any  ${\bf X}$  we want to estimate the expected value of the corresponding Y
- Suppose we want to predict the weight of a person from height

$$f^*(\mathbf{X}) = \mathbb{E}(Y|\mathbf{X})$$

- $\bullet$  Hence, for any  ${\bf X}$  we want to estimate the expected value of the corresponding Y
- Suppose we want to predict the weight of a person from height
- In reality, there are many people with the same height but different weight

$$f^*(\mathbf{X}) = \mathbb{E}(Y|\mathbf{X})$$

- $\bullet\,$  Hence, for any  ${\bf X}$  we want to estimate the expected value of the corresponding Y
- Suppose we want to predict the weight of a person from height
- In reality, there are many people with the same height but different weight
- Height does not determine the weight uniquely

$$f^*(\mathbf{X}) = \mathbb{E}(Y|\mathbf{X})$$

- $\bullet\,$  Hence, for any  ${\bf X}$  we want to estimate the expected value of the corresponding Y
- Suppose we want to predict the weight of a person from height
- In reality, there are many people with the same height but different weight
- Height does not determine the weight uniquely
- Relationship between height and weight is non-deterministic

$$f^*(\mathbf{X}) = \mathbb{E}(Y|\mathbf{X})$$

- $\bullet\,$  Hence, for any  ${\bf X}$  we want to estimate the expected value of the corresponding Y
- Suppose we want to predict the weight of a person from height
- In reality, there are many people with the same height but different weight
- Height does not determine the weight uniquely
- Relationship between height and weight is non-deterministic
- In case of mean regression, we predict the expected weight of people with a given height



2

• • • • • • • • • • •



< ∃⇒

э

• • • • • • • • • •
• Median regression:

Example: for a given height, predict the weight such that approximately half of the people with this height would be heavier than this

• Median regression:

Example: for a given height, predict the weight such that approximately half of the people with this height would be heavier than this

#### • Quantile regression:

Median regression generalized to any other quantile (median is the  $50\%\mbox{-}{\rm quantile}\mbox{)}$ 

• How do we evaluate how well  $\hat{f}$  approximates  $f^*(\mathbf{X}) = \mathbb{E}(Y|\mathbf{X})$ ?

# Evaluation of regression

- How do we evaluate how well  $\hat{f}$  approximates  $f^*(\mathbf{X}) = \mathbb{E}(Y|\mathbf{X})$ ?
- Even if we have hold-out test data, we still do not know the true  $f^*(\mathbf{X}) = \mathbb{E}(Y|\mathbf{X}).$

# Evaluation of regression

- How do we evaluate how well  $\hat{f}$  approximates  $f^*(\mathbf{X}) = \mathbb{E}(Y|\mathbf{X})$ ?
- Even if we have hold-out test data, we still do not know the true  $f^*(\mathbf{X}) = \mathbb{E}(Y|\mathbf{X}).$
- $f^*(\mathbf{X})$  minimizes the expected squared error on future data

$$f^*(\mathbf{X}) = \operatorname*{argmin}_{y \in \mathbb{R}} \mathbb{E}[(y - Y)^2 | \mathbf{X}] = \mathbb{E}[Y | \mathbf{X}]$$

# Evaluation of regression

- How do we evaluate how well  $\hat{f}$  approximates  $f^*(\mathbf{X}) = \mathbb{E}(Y|\mathbf{X})$ ?
- Even if we have hold-out test data, we still do not know the true  $f^*(\mathbf{X}) = \mathbb{E}(Y|\mathbf{X}).$
- $f^*(\mathbf{X})$  minimizes the expected squared error on future data

$$f^*(\mathbf{X}) = \operatorname*{argmin}_{y \in \mathbb{R}} \mathbb{E}[(y - Y)^2 | \mathbf{X}] = \mathbb{E}[Y | \mathbf{X}]$$

 Therefore, the most usual evaluation measure in (mean) regression is mean squared error (MSE) on test data:

$$\frac{1}{|Te|} \sum_{(\mathbf{x},y)\in Te} (\hat{f}(\mathbf{x}) - y)^2$$

• Noise is the difference between the true label and the prediction  $f^* = \mathbb{E}[Y|\mathbf{X}]$ 

$$\epsilon = Y - f^*(\mathbf{X})$$

< 1 k

э

• Noise is the difference between the true label and the prediction  $f^* = \mathbb{E}[Y|\mathbf{X}]$ 

$$\epsilon = Y - f^*(\mathbf{X})$$

• This relationship is usually presented as:

$$Y = f^*(X) + \epsilon$$

• Noise is the difference between the true label and the prediction  $f^* = \mathbb{E}[Y|\mathbf{X}]$ 

$$\epsilon = Y - f^*(\mathbf{X})$$

• This relationship is usually presented as:

$$Y = f^*(X) + \epsilon$$

• The task in mean regression is to predict  $f^*(X)$  and the noise cannot be (and should not attempted to be) predicted from the features X

• Suppose we want to predict a person's weight from height

< 行

э

- Suppose we want to predict a person's weight from height
- What is the noise here?

- Suppose we want to predict a person's weight from height
- What is the noise here?
- Noise is the difference between the actual weight and the average weight among people with the same height

### More on evaluation of regression

• Two common measures:

э

#### More on evaluation of regression

- Two common measures:
  - MSE mean squared error

$$MSE = \frac{1}{|Te|} \sum_{(\mathbf{x}, y) \in Te} (\hat{f}(\mathbf{x}) - y)^2$$

### More on evaluation of regression

- Two common measures:
  - MSE mean squared error

$$MSE = \frac{1}{|Te|} \sum_{(\mathbf{x}, y) \in Te} (\hat{f}(\mathbf{x}) - y)^2$$

• RMSE - root mean squared error

$$RMSE = \sqrt{MSE}$$

- Two common measures:
  - MSE mean squared error

$$MSE = \frac{1}{|Te|} \sum_{(\mathbf{x}, y) \in Te} (\hat{f}(\mathbf{x}) - y)^2$$

• RMSE - root mean squared error

$$RMSE = \sqrt{MSE}$$

• The advantage of RMSE: easier to interpret because it is measured in the same units as the target variable (whereas MSE is measured in these units squared)

• In regression the task is to learn a function approximator:  $\hat{f}:\mathbb{X}\to\mathbb{R}$ 

- $\bullet$  In regression the task is to learn a function approximator:  $\hat{f}:\mathbb{X}\to\mathbb{R}$
- In linear regression:

- In regression the task is to learn a function approximator:  $\hat{f}:\mathbb{X}\to\mathbb{R}$
- In linear regression:
  - We assume that the features are all numeric:

- In regression the task is to learn a function approximator:  $\hat{f}:\mathbb{X}\to\mathbb{R}$
- In linear regression:
  - We assume that the features are all numeric:

• We must learn a linear model:

$$\hat{f}(\mathbf{x}) = w_0 + \mathbf{w} \cdot \mathbf{x} = w_0 + \sum_{i=1}^d w_i x_i$$

- In regression the task is to learn a function approximator:  $\hat{f}:\mathbb{X}\to\mathbb{R}$
- In linear regression:
  - We assume that the features are all numeric:

• We must learn a linear model:

$$\hat{f}(\mathbf{x}) = w_0 + \mathbf{w} \cdot \mathbf{x} = w_0 + \sum_{i=1}^d w_i x_i$$

•  $w_0$  - intercept

- In regression the task is to learn a function approximator:  $\hat{f}:\mathbb{X}\to\mathbb{R}$
- In linear regression:
  - We assume that the features are all numeric:

• We must learn a linear model:

$$\hat{f}(\mathbf{x}) = w_0 + \mathbf{w} \cdot \mathbf{x} = w_0 + \sum_{i=1}^d w_i x_i$$

w<sub>0</sub> - intercept
w - coefficients (coefficient vector)

• Linear regression with a single feature

- Linear regression with a single feature
- We must learn a univariate linear function:

$$\hat{f}(\mathbf{x}) = w_0 + \mathbf{w} \cdot \mathbf{x} = w_0 + w_1 x_1,$$

- Linear regression with a single feature
- We must learn a univariate linear function:

$$\hat{f}(\mathbf{x}) = w_0 + \mathbf{w} \cdot \mathbf{x} = w_0 + w_1 x_1,$$

• We have training data:  $(x_1,y_1), (x_2,y_2), \ldots, (x_n,y_n) \in \mathbb{R}^2$ 

- Linear regression with a single feature
- We must learn a univariate linear function:

$$\hat{f}(\mathbf{x}) = w_0 + \mathbf{w} \cdot \mathbf{x} = w_0 + w_1 x_1,$$

- We have training data:  $(x_1,y_1),(x_2,y_2),\ldots,(x_n,y_n)\in\mathbb{R}^2$
- Task is to learn  $w_0$  and  $w_1$  such that  $\hat{f}$  minimizes future squared error

- Linear regression with a single feature
- We must learn a univariate linear function:

$$\hat{f}(\mathbf{x}) = w_0 + \mathbf{w} \cdot \mathbf{x} = w_0 + w_1 x_1,$$

- We have training data:  $(x_1,y_1),(x_2,y_2),\ldots,(x_n,y_n)\in\mathbb{R}^2$
- Task is to learn  $w_0$  and  $w_1$  such that  $\hat{f}$  minimizes future squared error
- One of the popular methods is the Ordinary Least Squares (OLS) method

# Ordinary Least Squares (OLS)



# Univariate OLS

Ordinary least squares (OLS) regression learns the weights by minimizing MSE on training data:

$$\hat{w}_0, \hat{w}_1 = \operatorname*{argmin}_{w_0, w_1} MSE(w_0, w_1) = \frac{1}{|Tr|} \sum_{(x,y)\in Tr} ((w_0 + w_1x) - y)^2$$



$$MSE(w_0, w_1) = \frac{1}{|Tr|} \sum_{(x,y)\in Tr} ((w_0 + w_1x) - y)^2$$

• Calculating the gradient with respect to  $w_0$ 

$$\frac{\partial MSE}{\partial w_0} =$$

FAST Foundation

$$MSE(w_0, w_1) = \frac{1}{|Tr|} \sum_{(x,y)\in Tr} ((w_0 + w_1x) - y)^2$$

• Calculating the gradient with respect to  $w_0$ 

$$\frac{\partial MSE}{\partial w_0} = \frac{1}{|Tr|} \sum_{(x,y)\in Tr} 2(w_0 + w_1 x - y) = 0$$

FAST Foundation

$$MSE(w_0, w_1) = \frac{1}{|Tr|} \sum_{(x,y)\in Tr} ((w_0 + w_1x) - y)^2$$

• Calculating the gradient with respect to  $w_0$ 

$$\frac{\partial MSE}{\partial w_0} = \frac{1}{|Tr|} \sum_{(x,y)\in Tr} 2(w_0 + w_1x - y) = 0$$
$$\frac{1}{|Tr|} \sum_{(x,y)\in Tr} y = w_0 + w_1 \frac{1}{|Tr|} \sum_{(x,y)\in Tr} x$$

# Calculating the partial derivatives

$$MSE(w_0, w_1) = \frac{1}{|Tr|} \sum_{(x,y)\in Tr} ((w_0 + w_1x) - y)^2$$

• Calculating the gradient with respect to  $w_0$ 

$$\begin{split} \frac{\partial MSE}{\partial w_0} &= \frac{1}{|Tr|} \sum_{(x,y)\in Tr} 2(w_0 + w_1x - y) = 0\\ \frac{1}{|Tr|} \sum_{(x,y)\in Tr} y &= w_0 + w_1 \frac{1}{|Tr|} \sum_{(x,y)\in Tr} x\\ w_0 &= \bar{y} - w_1 \bar{x} \end{split}$$
  
where  $\bar{x} &= \frac{1}{|Tr|} \sum_{(x,y)\in Tr} x$  and  $\bar{y} = \frac{1}{|Tr|} \sum_{(x,y)\in Tr} y$ 

$$MSE(w_0, w_1) = \frac{1}{|Tr|} \sum_{(x,y) \in Tr} ((w_0 + w_1x) - y)^2$$
  
Calculating the gradient with respect to  $w_1$ 

 $\frac{\partial MSE}{\partial w_1} =$ 

Image: A matrix and a matrix

æ

$$MSE(w_0, w_1) = \frac{1}{|Tr|} \sum_{(x,y)\in Tr} ((w_0 + w_1 x) - y)^2$$

Calculating the gradient with respect to  $w_1$ 

$$\frac{\partial MSE}{\partial w_1} = \frac{1}{|Tr|} \sum_{(x,y)\in Tr} 2(w_0 + w_1x - y)x = 0$$

$$MSE(w_0, w_1) = \frac{1}{|Tr|} \sum_{(x,y)\in Tr} ((w_0 + w_1x) - y)^2$$

Calculating the gradient with respect to  $w_1$ 

$$\frac{\partial MSE}{\partial w_1} = \frac{1}{|Tr|} \sum_{(x,y)\in Tr} 2(w_0 + w_1x - y)x = 0$$
$$\frac{1}{|Tr|} \sum_{(x,y)\in Tr} (\bar{y} - w_1\bar{x} + w_1x - y)x = 0$$
$$MSE(w_0, w_1) = \frac{1}{|Tr|} \sum_{(x,y)\in Tr} ((w_0 + w_1x) - y)^2$$

Calculating the gradient with respect to  $w_1$ 

$$\frac{\partial MSE}{\partial w_1} = \frac{1}{|Tr|} \sum_{(x,y)\in Tr} 2(w_0 + w_1x - y)x = 0$$
$$\frac{1}{|Tr|} \sum_{(x,y)\in Tr} (\bar{y} - w_1\bar{x} + w_1x - y)x = 0$$
$$w_1 \sum_{(x,y)\in Tr} (x - \bar{x}) \cdot x = \sum_{(x,y)\in Tr} (y - \bar{y}) \cdot x$$

$$MSE(w_0, w_1) = \frac{1}{|Tr|} \sum_{(x,y)\in Tr} ((w_0 + w_1x) - y)^2$$

Calculating the gradient with respect to  $w_1$ 

$$\frac{\partial MSE}{\partial w_1} = \frac{1}{|Tr|} \sum_{(x,y)\in Tr} 2(w_0 + w_1x - y)x = 0$$
$$\frac{1}{|Tr|} \sum_{(x,y)\in Tr} (\bar{y} - w_1\bar{x} + w_1x - y)x = 0$$
$$w_1 \sum_{(x,y)\in Tr} (x - \bar{x}) \cdot x = \sum_{(x,y)\in Tr} (y - \bar{y}) \cdot x$$
$$\left[\sum(x - \bar{x}) \cdot x - \sum(x - \bar{x}) \cdot \bar{x}\right] = \sum(y - \bar{y}) \cdot x - \sum(y - \bar{y}) \cdot \bar{x}$$

FAST Foundation

 $w_1($ 

э

## Calculating the partial derivatives

$$MSE(w_0, w_1) = \frac{1}{|Tr|} \sum_{(x,y)\in Tr} ((w_0 + w_1 x) - y)^2$$

Calculating the gradient with respect to  $w_1$ 

$$\frac{\partial MSE}{\partial w_1} = \frac{1}{|Tr|} \sum_{(x,y)\in Tr} 2(w_0 + w_1x - y)x = 0$$
$$\frac{1}{|Tr|} \sum_{(x,y)\in Tr} (\bar{y} - w_1\bar{x} + w_1x - y)x = 0$$
$$w_1 \sum_{(x,y)\in Tr} (x - \bar{x}) \cdot x = \sum_{(x,y)\in Tr} (y - \bar{y}) \cdot x$$
$$w_1 \Big( \sum (x - \bar{x}) \cdot x - \sum (x - \bar{x}) \cdot \bar{x} \Big) = \sum (y - \bar{y}) \cdot x - \sum (y - \bar{y}) \cdot \bar{x}$$
$$w_1 \sum (x - \bar{x})^2 = \sum (y - \bar{y})(x - \bar{x}) \Rightarrow$$

**FAST** Foundation

 $w_1$ 

## Calculating the partial derivatives

$$MSE(w_0, w_1) = \frac{1}{|Tr|} \sum_{(x,y)\in Tr} ((w_0 + w_1 x) - y)^2$$

Calculating the gradient with respect to  $w_1$ 

$$\begin{aligned} \frac{\partial MSE}{\partial w_1} &= \frac{1}{|Tr|} \sum_{(x,y)\in Tr} 2(w_0 + w_1x - y)x = 0\\ \frac{1}{|Tr|} \sum_{(x,y)\in Tr} (\bar{y} - w_1\bar{x} + w_1x - y)x = 0\\ w_1 \sum_{(x,y)\in Tr} (x - \bar{x}) \cdot x &= \sum_{(x,y)\in Tr} (y - \bar{y}) \cdot x\\ w_1 \Big( \sum (x - \bar{x}) \cdot x - \sum (x - \bar{x}) \cdot \bar{x} \Big) &= \sum (y - \bar{y}) \cdot x - \sum (y - \bar{y}) \cdot \bar{x}\\ w_1 \sum (x - \bar{x})^2 &= \sum (y - \bar{y})(x - \bar{x}) \Rightarrow \quad w_1 = \frac{\sum (y - \bar{y})(x - \bar{x})}{\sum (x - \bar{x})^2} = \end{aligned}$$

### Calculating the partial derivatives

$$MSE(w_0, w_1) = \frac{1}{|Tr|} \sum_{(x,y)\in Tr} ((w_0 + w_1 x) - y)^2$$

Calculating the gradient with respect to  $w_1$ 

$$\frac{\partial MSE}{\partial w_1} = \frac{1}{|Tr|} \sum_{(x,y)\in Tr} 2(w_0 + w_1x - y)x = 0$$
$$\frac{1}{|Tr|} \sum_{(x,y)\in Tr} (\bar{y} - w_1\bar{x} + w_1x - y)x = 0$$
$$w_1 \sum_{(x,y)\in Tr} (x - \bar{x}) \cdot x = \sum_{(x,y)\in Tr} (y - \bar{y}) \cdot x$$
$$w_1 \Big( \sum (x - \bar{x}) \cdot x - \sum (x - \bar{x}) \cdot \bar{x} \Big) = \sum (y - \bar{y}) \cdot x - \sum (y - \bar{y}) \cdot \bar{x}$$
$$w_1 \sum (x - \bar{x})^2 = \sum (y - \bar{y}) \cdot x - \sum (y - \bar{y}) \cdot \bar{x} = \sum (y - \bar{y}) \cdot x - \sum (y - \bar{y}) \cdot \bar{x}$$

# Univariate OLS

Ordinary least squares (OLS) regression learns the weights by minimizing MSE on training data:

$$\hat{w}_{0}, \hat{w}_{1} = \operatorname*{argmin}_{w_{0}, w_{1}} MSE(w_{0}, w_{1}) = \begin{cases} \hat{w}_{1} &= \frac{Cov(x, y)}{Var(x)} \\ \hat{w}_{0} &= \bar{y} - \hat{w}_{1}\bar{x} \end{cases}$$



• Fitting: Calculate  $\hat{w}_0, \hat{w}_1$  based on the

э

Image: A matrix and a matrix

Calculate  $\hat{w}_0, \hat{w}_1$  based on the

 ${\ensuremath{\,\circ}}$  sample mean and variance of feature values x

< 1 k

Calculate  $\hat{w}_0, \hat{w}_1$  based on the

- ${\ensuremath{\, \bullet }}$  sample mean and variance of feature values x
- ${\ensuremath{\, \rm o}}$  sample mean of y

Calculate  $\hat{w}_0, \hat{w}_1$  based on the

- ${\ensuremath{\, \bullet }}$  sample mean and variance of feature values x
- ${\scriptstyle \bullet} \,$  sample mean of y
- $\bullet\,$  sample covariance of x and y

Calculate  $\hat{w}_0, \hat{w}_1$  based on the

- ${\ensuremath{\, \bullet }}$  sample mean and variance of feature values x
- ${\scriptstyle \bullet} \,$  sample mean of y
- $\bullet\,$  sample covariance of x and y
- Predicting for a new instance *x*:

$$\hat{y} = \hat{f}(x) = \hat{w}_0 + \hat{w}_1 x$$

Calculate  $\hat{w}_0, \hat{w}_1$  based on the

- ${\ensuremath{\, \bullet }}$  sample mean and variance of feature values x
- ${\scriptstyle \bullet} \,$  sample mean of y
- ${\ensuremath{\, \circ }}$  sample covariance of x and y
- Predicting for a new instance *x*:

$$\hat{y} = \hat{f}(x) = \hat{w}_0 + \hat{w}_1 x$$

• If before learning the regression model we **standardise** both the feature and the target variable (zero mean and unit variance), then

# Univariate OLS

• Fitting:

Calculate  $\hat{w}_0, \hat{w}_1$  based on the

- $\bullet\,$  sample mean and variance of feature values x
- sample mean of y
- sample covariance of  $\boldsymbol{x}$  and  $\boldsymbol{y}$
- Predicting for a new instance *x*:

$$\hat{y} = \hat{f}(x) = \hat{w}_0 + \hat{w}_1 x$$

• If before learning the regression model we **standardise** both the feature and the target variable (zero mean and unit variance), then

$$\hat{w}_0, \hat{w}_1 = \operatorname*{argmin}_{w_0, w_1} MSE(w_0, w_1) = \begin{cases} \hat{w}_1 = \frac{\hat{Cov}(x, y)}{\hat{Var}(x)} = \hat{Corr}(x, y) \\ \hat{w}_0 = \bar{y} - w_1 \bar{x} = 0 \end{cases}$$

# Ordinary Least Squares

- OLS is very sensitive to outliers
- A single faraway point can significantly shift the predictions



# Ordinary Least Squares

- OLS is very sensitive to outliers
- A single faraway point can significantly shift the predictions



### Multivariate OLS

• Linear model:

$$\hat{f}(\mathbf{x}) = w_0 + \mathbf{w} \cdot \mathbf{x} = w_0 + \sum_{i=1}^d w_i x_i$$

イロト イヨト イヨト イヨト

2

• Linear model:

$$\hat{f}(\mathbf{x}) = w_0 + \mathbf{w} \cdot \mathbf{x} = w_0 + \sum_{i=1}^d w_i x_i$$

• Ordinary least squares (OLS) regression learns the weights by minimizing MSE on training data:

$$\hat{w}_0, \hat{\mathbf{w}} = \operatorname*{argmin}_{w_0, \mathbf{w}} MSE(w_0, \mathbf{w}) = \operatorname*{argmin}_{w_0, \mathbf{w}} \frac{1}{|Tr|} \sum_{(x, y) \in Tr} ((w_0 + \mathbf{w} \cdot \mathbf{x}) - y)^2$$

• Linear model:

$$\hat{f}(\mathbf{x}) = w_0 + \mathbf{w} \cdot \mathbf{x} = w_0 + \sum_{i=1}^d w_i x_i$$

• Ordinary least squares (OLS) regression learns the weights by minimizing MSE on training data:

$$\hat{w}_0, \hat{\mathbf{w}} = \operatorname*{argmin}_{w_0, \mathbf{w}} MSE(w_0, \mathbf{w}) = \operatorname*{argmin}_{w_0, \mathbf{w}} \frac{1}{|Tr|} \sum_{(x, y) \in Tr} ((w_0 + \mathbf{w} \cdot \mathbf{x}) - y)^2$$

• Let's add a feature, which is always 1

$$\mathbf{x} = (x_1, \dots, x_d) \to \mathbf{x'} = (1, x_1, \dots, x_d) = (x_0, x_1, \dots, x_d)$$

## Multivariate OLS

• Linear model:

$$\hat{f}(\mathbf{x}) = w_0 + \mathbf{w} \cdot \mathbf{x} = w_0 + \sum_{i=1}^d w_i x_i$$

• Ordinary least squares (OLS) regression learns the weights by minimizing MSE on training data:

$$\hat{w}_0, \hat{\mathbf{w}} = \operatorname*{argmin}_{w_0, \mathbf{w}} MSE(w_0, \mathbf{w}) = \operatorname*{argmin}_{w_0, \mathbf{w}} \frac{1}{|Tr|} \sum_{(x, y) \in Tr} ((w_0 + \mathbf{w} \cdot \mathbf{x}) - y)^2$$

• Let's add a feature, which is always 1

$$\mathbf{x} = (x_1, \dots, x_d) \to \mathbf{x}' = (1, x_1, \dots, x_d) = (x_0, x_1, \dots, x_d)$$

$$\hat{f}(\mathbf{x}') = \mathbf{w}' \cdot \mathbf{x}' = \sum_{i=0}^{d} w_i x_i$$

## Multivariate OLS

• Linear model:

$$\hat{f}(\mathbf{x}) = w_0 + \mathbf{w} \cdot \mathbf{x} = w_0 + \sum_{i=1}^d w_i x_i$$

• Ordinary least squares (OLS) regression learns the weights by minimizing MSE on training data:

$$\hat{w}_0, \hat{\mathbf{w}} = \operatorname*{argmin}_{w_0, \mathbf{w}} MSE(w_0, \mathbf{w}) = \operatorname*{argmin}_{w_0, \mathbf{w}} \frac{1}{|Tr|} \sum_{(x, y) \in Tr} ((w_0 + \mathbf{w} \cdot \mathbf{x}) - y)^2$$

• Let's add a feature, which is always 1

$$\mathbf{x} = (x_1, \dots, x_d) \to \mathbf{x}' = (1, x_1, \dots, x_d) = (x_0, x_1, \dots, x_d)$$

$$\hat{f}(\mathbf{x}') = \mathbf{w}' \cdot \mathbf{x}' = \sum_{i=0}^{d} w_i x_i$$

from now on we will skip the ' for simplicity

FAST Foundation

$$\hat{\mathbf{w}} = \operatorname*{argmin}_{\mathbf{w}} MSE(\mathbf{w}) = \operatorname*{argmin}_{\mathbf{w}} \frac{1}{|Tr|} \sum_{(\mathbf{x}, y) \in Tr} (\mathbf{w} \cdot \mathbf{x} - y)^2$$

2

メロト メポト メヨト メヨト

$$\hat{\mathbf{w}} = \operatorname*{argmin}_{\mathbf{w}} MSE(\mathbf{w}) = \operatorname*{argmin}_{\mathbf{w}} \frac{1}{|Tr|} \sum_{(\mathbf{x}, y) \in Tr} (\mathbf{w} \cdot \mathbf{x} - y)^2$$

• Denote the training instances by  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n) \in \mathbb{R}^{d+1} \times \mathbb{R}$ 

Image: A matrix and a matrix

э

$$\hat{\mathbf{w}} = \operatorname*{argmin}_{\mathbf{w}} MSE(\mathbf{w}) = \operatorname*{argmin}_{\mathbf{w}} \frac{1}{|Tr|} \sum_{(\mathbf{x}, y) \in Tr} (\mathbf{w} \cdot \mathbf{x} - y)^2$$

- Denote the training instances by  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n) \in \mathbb{R}^{d+1} \times \mathbb{R}$
- Denote the residual errors on instances by:

$$e_1 = y_1 - \mathbf{x}_1 \cdot \mathbf{w}, \dots, e_n = y_n - \mathbf{x}_n \cdot \mathbf{w}$$

$$\hat{\mathbf{w}} = \operatorname*{argmin}_{\mathbf{w}} MSE(\mathbf{w}) = \operatorname*{argmin}_{\mathbf{w}} \frac{1}{|Tr|} \sum_{(\mathbf{x}, y) \in Tr} (\mathbf{w} \cdot \mathbf{x} - y)^2$$

- Denote the training instances by  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n) \in \mathbb{R}^{d+1} \times \mathbb{R}$
- Denote the residual errors on instances by:

$$e_1 = y_1 - \mathbf{x}_1 \cdot \mathbf{w}, \dots, e_n = y_n - \mathbf{x}_n \cdot \mathbf{w}$$

 $\bullet\,$  In matrix form:  $\mathbf{e}=\mathbf{y}-\mathbf{X}\mathbf{w}$ 

$$\hat{\mathbf{w}} = \operatorname*{argmin}_{\mathbf{w}} \sum_{i=1}^{n} e_i^2 = \operatorname*{argmin}_{\mathbf{w}} \mathbf{e} \cdot \mathbf{e} = \operatorname*{argmin}_{\mathbf{w}} (\mathbf{y} - \mathbf{X}\mathbf{w}) \cdot (\mathbf{y} - \mathbf{X}\mathbf{w})$$

$$\hat{\mathbf{w}} = \operatorname*{argmin}_{\mathbf{w}} MSE(\mathbf{w}) = \operatorname*{argmin}_{\mathbf{w}} \frac{1}{|Tr|} \sum_{(\mathbf{x}, y) \in Tr} (\mathbf{w} \cdot \mathbf{x} - y)^2$$

- Denote the training instances by  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n) \in \mathbb{R}^{d+1} \times \mathbb{R}$
- Denote the residual errors on instances by:

$$e_1 = y_1 - \mathbf{x}_1 \cdot \mathbf{w}, \dots, e_n = y_n - \mathbf{x}_n \cdot \mathbf{w}$$

• In matrix form:  $\mathbf{e} = \mathbf{y} - \mathbf{X}\mathbf{w}$ 

$$\hat{\mathbf{w}} = \operatorname*{argmin}_{\mathbf{w}} \sum_{i=1}^{n} e_i^2 = \operatorname*{argmin}_{\mathbf{w}} \mathbf{e} \cdot \mathbf{e} = \operatorname*{argmin}_{\mathbf{w}} (\mathbf{y} - \mathbf{X}\mathbf{w}) \cdot (\mathbf{y} - \mathbf{X}\mathbf{w})$$

To solve this we equate the gradient to zero:

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

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

- ∢ 🗗 ▶

э

$$\frac{\partial (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})}{\partial \mathbf{w}} = -2\mathbf{X}^T (\mathbf{y} - \mathbf{X}\mathbf{w}) = 0$$

$$\frac{\partial (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})}{\partial \mathbf{w}} = -2\mathbf{X}^T (\mathbf{y} - \mathbf{X}\mathbf{w}) = 0$$
$$\mathbf{X}^T \mathbf{X}\mathbf{w} = \mathbf{X}^T \mathbf{y}$$

$$\frac{\partial (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})}{\partial \mathbf{w}} = -2\mathbf{X}^T (\mathbf{y} - \mathbf{X}\mathbf{w}) = 0$$
$$\mathbf{X}^T \mathbf{X}\mathbf{w} = \mathbf{X}^T \mathbf{y}$$
$$(\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X}) \mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{y})$$

$$\frac{\partial (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})}{\partial \mathbf{w}} = -2\mathbf{X}^T (\mathbf{y} - \mathbf{X}\mathbf{w}) = 0$$
$$\mathbf{X}^T \mathbf{X}\mathbf{w} = \mathbf{X}^T \mathbf{y}$$
$$(\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X}) \mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{y})$$
$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{y})$$

$$\frac{\partial (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})}{\partial \mathbf{w}} = -2\mathbf{X}^T (\mathbf{y} - \mathbf{X}\mathbf{w}) = 0$$
$$\mathbf{X}^T \mathbf{X}\mathbf{w} = \mathbf{X}^T \mathbf{y}$$
$$(\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X}) \mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{y})$$
$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{y})$$

For multivariate OLS there exists a **closed-form solution** (i.e. can be explicitly calculated without numerical optimization):

$$\hat{\mathbf{w}} = \operatorname*{argmin}_{\mathbf{w}} (\mathbf{y} - \mathbf{X}\mathbf{w}) \cdot (\mathbf{y} - \mathbf{X}\mathbf{w}) = (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{y})$$

**Pro:** if the number of instances is much bigger than the number of features, then OLS works quite well

- **Pro:** if the number of instances is much bigger than the number of features, then OLS works quite well
- **Con:** otherwise OLS tends to overfit the noise, particularly if there is a lot of noise in the data

- **Pro:** if the number of instances is much bigger than the number of features, then OLS works quite well
- **Con:** otherwise OLS tends to overfit the noise, particularly if there is a lot of noise in the data
- **Con:** if many features are collinear (highly correlated) then OLS tends to overfit

- $\checkmark$  Main concepts in regression
- ✓ Linear Regression
- ✓ Ordinary Least Squares (OLS)