

# Machine Learning for Networks: Regression

**Andrea Araldo**

*January 4, 2022*



	ML task	Linear Regression	Logistic Regression	Tree-based learning	Neural Networks	$k$ -Nearest Neighbors
Supervised	Regression Classification	x	x	x	x	
Unsupervised	Clustering Dimensionality reduction Anomaly detection			x	x	x
	Recommender Systems				x	

- Supervised Learning Models
- Linear Regression
- Train / Test / Cross-Validation

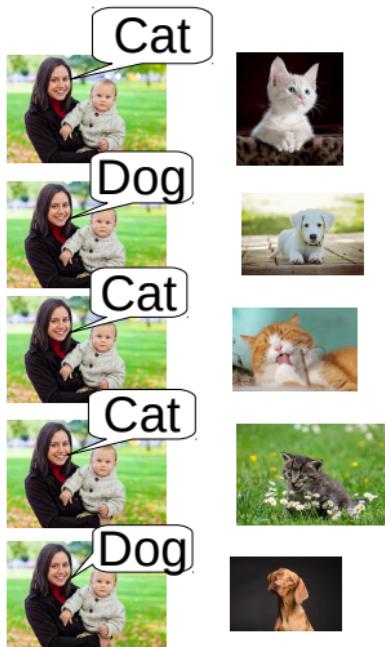
## Section 1

# **Introduction to Supervised Learning**

# How humans learn

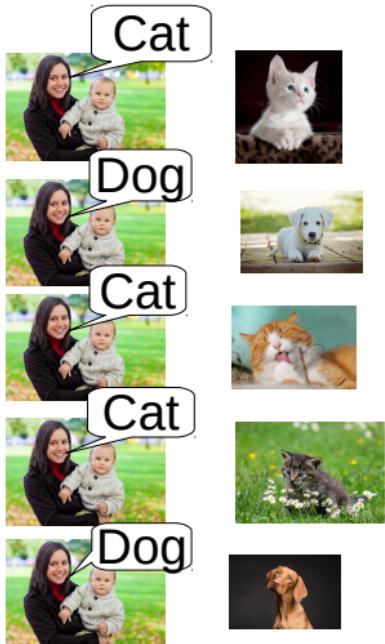
4 / 54

TRAINING (learning)



Some think instead that we learn in a different [Tho20]

## TRAINING (learning)



## INFERENCE (prediction)



Some think instead that we learn in a different [Tho20]

## TRAINING (learning)



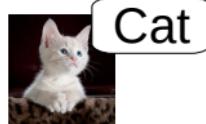
## INFERENCE (prediction)



Some think instead that we learn in a different [Tho20]

## TRAINING (learning)

Training dataset



## TRAINING (learning)

### Training dataset

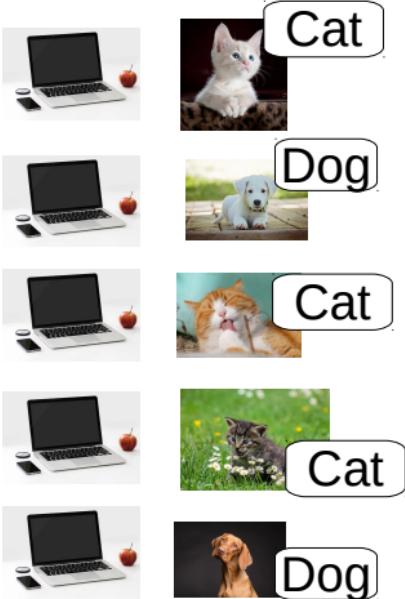


### INFERENCE (prediction)



## TRAINING (learning)

### Training dataset

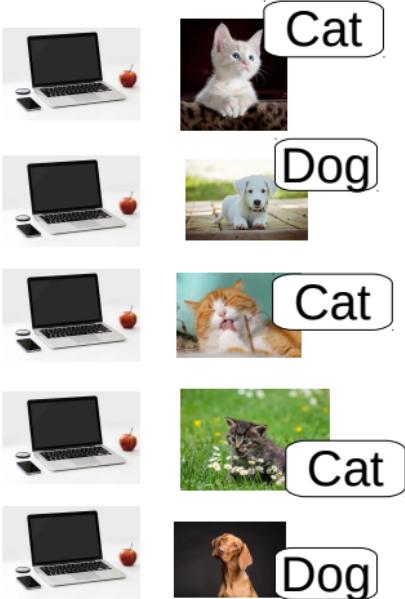


### INFERENCE (prediction)



## TRAINING (learning)

### Training dataset



### INFERENCE (prediction)



Deep Neural Network better than humans to classify several image datasets.

Cambridge English Dictionary:

*the process of getting an understanding of something by studying it or by experience*

Cambridge English Dictionary:

*the process of getting an understanding of something by studying it or by experience*

What can we learn in a network, by just observing traffic (encrypted)?

- Anomaly/attack
- Application protocols hidden in encrypted traffic
- Quality of Experience (QoE)
- Malfunctioning of a device
- ...

- **Problem [KEB19]:**

*Infer* video chunk resolution

*only observing* network information, e.g., throughput, packet size, inter-arrival times, etc.

Traffic is encrypted.

- **Motivation:**

Network operators can have insights on QoE without *Deep Packet Inspection*

*Features (or Columns):*

- Throughput mean (bps)  
Throughput st.dev. (bps)  
Packet size mean (B)  
Packet size st.dev.  
Inter-arrival mean (ms), etc.

*Label (or target):*

- Resolution of downloaded video chunks (360p, 720p)
- (it indicates the number of vertical pixels)

*Features (or Columns):*

- Throughput mean (bps)  
Throughput st.dev. (bps)  
Packet size mean (B)  
Packet size st.dev.  
Inter-arrival mean (ms), etc.

*Label (or target):*

- Resolution of downloaded video chunks (360p, 720p)
- (it indicates the number of vertical pixels)

*Model:*

- A function  $h(\text{features}) = \text{label}$
- If label  $\in \mathbb{R}$ : **regression** model
- If label discrete: **classification** model

# Example

9 / 54

		Features (or Independent variables)					True Label	
	Sample	TP mean	TP st.dev.	Pkt size mean	Pkt std. dev.	Int-Arriv mean	$y^{(i)}$ (Resolution)	
Our dataset	$\mathbf{x}^{(1)^T} =$ $\vdots$ $\mathbf{x}^{(M)^T} =$	2 Mbps ⋮ 2 Mbps	1 Kbps ⋮ 1.1 Kbps	1 KB ⋮ 1.3 KB	300 B ⋮ 400 B	30 ms ⋮ 20 ms	<b>475</b> ⋮ <b>720</b>	
New samples	$\mathbf{x}^{(M+1)^T} =$ $\vdots$	2 Mbps ⋮	1 Kbps ⋮	1 KB ⋮	300 B ⋮	30 ms ⋮	? ⋮	

# Example

9 / 54

		Features (or Independent variables)					True Label	Pred Label	
	Sample	TP mean	TP st.dev.	Pkt size mean	Pkt std. dev.	Int-Arriv mean	$y^{(i)}$ (Resolution)	$\hat{y}^{(i)} = h(\mathbf{x}^{(i)})$	
Our dataset	$\mathbf{x}^{(1)^T} =$ ⋮ $\mathbf{x}^{(M)^T} =$	2 Mbps	1 Kbps	1 KB	300 B	30 ms	<b>475</b>	482	⋮
	$\mathbf{x}^{(M+1)^T} =$ ⋮	2 Mbps	1.1 Kbps	1.3 KB	400 B	20 ms	<b>720</b>	693	⋮
New samples							?	1078	⋮

- Model: function  $h$  used to predict

# Example

9 / 54

		Features (or Independent variables)					True Label	Pred Label	Pred Label
	Sample	TP mean	TP st.dev.	Pkt size mean	Pkt std. dev.	Int-Arriv mean	$y^{(i)}$ (Resolution)	$\hat{y}^{(i)} = h(\mathbf{x}^{(i)})$	$\hat{y}^{(i)} = h(\mathbf{x}^{(i)})$
Our dataset	$\mathbf{x}^{(1)^T} =$ ⋮ $\mathbf{x}^{(M)^T} =$	2 Mbps	1 Kbps	1 KB	300 B	30 ms	<b>475</b>	482	1083
	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
New samples	$\mathbf{x}^{(M+1)^T} =$ ⋮	2 Mbps	1 Kbps	1 KB	300 B	30 ms	<b>720</b>	693	323
	⋮	⋮	⋮	⋮	⋮	⋮	?	1078	376

- Model: function  $h$  used to predict

*Which of the two models would you trust more to predict?*

# Example

9 / 54

		Features (or Independent variables)					True Label	Pred Label	Pred Label
	Sample	TP mean	TP st.dev.	Pkt size mean	Pkt std. dev.	Int-Arriv mean	$y^{(i)}$ (Resolution)	$\hat{y}^{(i)} = h(\mathbf{x}^{(i)})$	$\hat{y}^{(i)} = h(\mathbf{x}^{(i)})$
Our dataset	$\mathbf{x}^{(1)^T} =$ ⋮ $\mathbf{x}^{(M)^T} =$	2 Mbps	1 Kbps	1 KB	300 B	30 ms	<b>475</b>	482	1083
		⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
New samples	$\mathbf{x}^{(M+1)^T} =$ ⋮	2 Mbps	1 Kbps	1 KB	300 B	30 ms	<b>720</b>	693	323
		⋮	⋮	⋮	⋮	⋮	?	1078	376

- Model: function  $h$  used to predict

*Which of the two models would you trust more to predict?*

- $h$  is “good” if  $h(\mathbf{x}^{(i)}) \simeq y^{(i)}$  for the new samples.

# Example

9 / 54

	Sample	Features (or Independent variables)					True Label	Pred Label	Pred Label
Our dataset	$\mathbf{x}^{(1)T} =$ ⋮ $\mathbf{x}^{(M)T} =$	TP mean ⋮ 2 Mbps	TP st.dev. ⋮ 1 Kbps	Pkt size mean ⋮ 1 KB	Pkt std. dev. ⋮ 300 B	Int-Arriv mean ⋮ 30 ms	$y^{(i)} \text{ (Resolution)}$ ⋮ <b>475</b>	$\hat{y}^{(i)} = h(\mathbf{x}^{(i)})$ ⋮ 482	1083 ⋮ 323
	$\mathbf{x}^{(M+1)T} =$ ⋮	2 Mbps ⋮	1 Kbps ⋮	1 KB ⋮	300 B ⋮	30 ms ⋮	<b>720</b> ⋮ ?	693 ⋮ 1078	376 ⋮
New samples									

- Model: function  $h$  used to predict

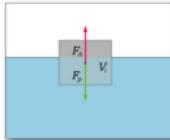
*Which of the two models would you trust more to predict?*

- $h$  is “good” if  $h(\mathbf{x}^{(i)}) \simeq y^{(i)}$  for the new samples.
- But we do not know  $y^{(i)}$  for the new samples.
- We can just evaluate how  $h(\cdot)$  performs in our dataset.

# Famous examples of models “by hand”

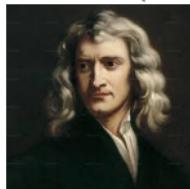
10 / 54

Archimedes (III cent. BC)



weight of displaced fluid = weight of object in vacuum – weight of object in fluid

Newton (XVII cent.)



$$\mathbf{F} = m\mathbf{a}$$

Maxwell (XVII cent.)



$$\nabla \cdot \mathbf{E} = \frac{\rho}{\epsilon_0}$$

$$\nabla \cdot \mathbf{B} = 0$$

$$\nabla \times \mathbf{E} = -\frac{\partial \mathbf{B}}{\partial t}$$

$$\nabla \times \mathbf{B} = \mu_0 \left( \mathbf{J} + \epsilon_0 \frac{\partial \mathbf{E}}{\partial t} \right)$$

They created a model  $h(\cdot)$  relating a target to some features, **based on observation.**

## Generalization

For **any** object with  $m$  and  $\mathbf{a}$ , we know the force  $\mathbf{F}$  it produces.

Pictures: Wikipedia

- Convention

$$\mathbf{x}^{(i)T} = (x_1^{(i)}, \dots, x_N^{(i)}) \in \mathbb{R}^{N+1}$$

- Convention

$$\mathbf{x}^{(i)T} = (1, x_1^{(i)}, \dots, x_N^{(i)}) \in \mathbb{R}^{N+1}$$

- Convention

$$\mathbf{x}^{(i)T} = (1, x_1^{(i)}, \dots, x_N^{(i)}) \in \mathbb{R}^{N+1}$$

- Linear Model

$$\begin{aligned} h_{\theta}(\mathbf{x}^{(i)}) &= \theta^T \cdot \mathbf{x}^{(i)} = \sum_{j=0}^N \theta_j \cdot x_j^{(i)} \\ &= \theta_0 + \theta_{\text{TP mean}} \cdot x_j^{(i)} \\ &\quad + \theta_{\text{TP st.dev.}} \cdot x_{\text{TP st.dev.}}^{(i)} + \dots \end{aligned}$$

- Convention

$$\mathbf{x}^{(i)T} = (1, x_1^{(i)}, \dots, x_N^{(i)}) \in \mathbb{R}^{N+1}$$

- Quadratic Model

$$h_{\theta}(\mathbf{x}^{(i)}) = \theta_0 + \theta_1 x_1^{(i)} + \dots + \theta_N x_N^{(i)}$$

- Linear Model

$$\begin{aligned} h_{\theta}(\mathbf{x}^{(i)}) &= \theta^T \cdot \mathbf{x}^{(i)} = \sum_{j=0}^N \theta_j \cdot x_j^{(i)} \\ &= \theta_0 + \theta_{\text{TP mean}} \cdot x_j^{(i)} \\ &\quad + \theta_{\text{TP st.dev.}} \cdot x_{\text{TP st.dev.}}^{(i)} + \dots \end{aligned}$$

- Convention

$$\mathbf{x}^{(i)T} = (1, x_1^{(i)}, \dots, x_N^{(i)}) \in \mathbb{R}^{N+1}$$

- Quadratic Model

$$\begin{aligned} h_{\theta}(\mathbf{x}^{(i)}) &= \theta_0 + \theta_1 x_1^{(i)} + \dots + \theta_N x_N^{(i)} \\ &\quad + \theta_{N+1} \cdot x_1^{(i)2} + \dots + \theta_N \cdot x_{N+N}^{(i)2} \end{aligned}$$

- Linear Model

$$\begin{aligned} h_{\theta}(\mathbf{x}^{(i)}) &= \theta^T \cdot \mathbf{x}^{(i)} = \sum_{j=0}^N \theta_j \cdot x_j^{(i)} \\ &= \theta_0 + \theta_{\text{TP mean}} \cdot x_j^{(i)} \\ &\quad + \theta_{\text{TP st.dev.}} \cdot x_{\text{TP st.dev.}}^{(i)} + \dots \end{aligned}$$

- Convention

$$\mathbf{x}^{(i)T} = (1, x_1^{(i)}, \dots, x_N^{(i)}) \in \mathbb{R}^{N+1}$$

- Linear Model

$$\begin{aligned} h_{\theta}(\mathbf{x}^{(i)}) &= \theta^T \cdot \mathbf{x}^{(i)} = \sum_{j=0}^N \theta_j \cdot x_j^{(i)} \\ &= \theta_0 + \theta_{\text{TP mean}} \cdot x_j^{(i)} \\ &\quad + \theta_{\text{TP st.dev.}} \cdot x_{\text{TP st.dev.}}^{(i)} + \dots \end{aligned}$$

- Quadratic Model

$$\begin{aligned} h_{\theta}(\mathbf{x}^{(i)}) &= \theta_0 + \theta_1 x_1^{(i)} + \dots + \theta_N x_N^{(i)} \\ &\quad + \theta_{N+1} \cdot x_1^{(i)2} + \dots + \theta_N \cdot x_{N+N}^{(i)2} \\ &\quad + \theta_{2N+1} \cdot x_1^{(i)} \cdot x_2^{(i)} + \theta_{2N+1} \cdot x_1^{(i)} \cdot x_{2N+1}^{(i)} \\ &\quad + \dots \end{aligned}$$

- Convention

$$\mathbf{x}^{(i)T} = (1, x_1^{(i)}, \dots, x_N^{(i)}) \in \mathbb{R}^{N+1}$$

- Linear Model

$$\begin{aligned} h_{\theta}(\mathbf{x}^{(i)}) &= \theta^T \cdot \mathbf{x}^{(i)} = \sum_{j=0}^N \theta_j \cdot x_j^{(i)} \\ &= \theta_0 + \theta_{\text{TP mean}} \cdot x_j^{(i)} \\ &\quad + \theta_{\text{TP st.dev.}} \cdot x_{\text{TP st.dev.}}^{(i)} + \dots \end{aligned}$$

- Quadratic Model

$$\begin{aligned} h_{\theta}(\mathbf{x}^{(i)}) &= \theta_0 + \theta_1 x_1^{(i)} + \dots + \theta_N x_N^{(i)} \\ &\quad + \theta_{N+1} \cdot x_1^{(i)2} + \dots + \theta_N \cdot x_{N+N}^{(i)2} \\ &\quad + \theta_{2N+1} \cdot x_1^{(i)} \cdot x_2^{(i)} + \theta_{2N+1} \cdot x_1^{(i)} \cdot x_{2N+1}^{(i)} \\ &\quad + \dots \end{aligned}$$

- Neural Network (NN):

$h_{\theta}(\cdot)$  is the output of a NN with weights  $\theta$ .

- Convention

$$\mathbf{x}^{(i)T} = (1, x_1^{(i)}, \dots, x_N^{(i)}) \in \mathbb{R}^{N+1}$$

- Linear Model

$$\begin{aligned} h_{\theta}(\mathbf{x}^{(i)}) &= \theta^T \cdot \mathbf{x}^{(i)} = \sum_{j=0}^N \theta_j \cdot x_j^{(i)} \\ &= \theta_0 + \theta_{\text{TP mean}} \cdot x_j^{(i)} \\ &\quad + \theta_{\text{TP st.dev.}} \cdot x_{\text{TP st.dev.}}^{(i)} + \dots \end{aligned}$$

**Loss:**  $J = \frac{1}{M} \sum_{i=1}^M \left( \underbrace{y^{(i)} - h_{\theta}(\mathbf{x}^{(i)})}_{\text{Residual } \varepsilon^{(i)}} \right)^2$

- Quadratic Model

$$\begin{aligned} h_{\theta}(\mathbf{x}^{(i)}) &= \theta_0 + \theta_1 x_1^{(i)} + \dots + \theta_N x_N^{(i)} \\ &\quad + \theta_{N+1} \cdot x_1^{(i)2} + \dots + \theta_N \cdot x_{N+N}^{(i)2} \\ &\quad + \theta_{2N+1} \cdot x_1^{(i)} \cdot x_2^{(i)} + \theta_{2N+1} \cdot x_1^{(i)} \cdot x_{2N+1}^{(i)} \\ &\quad + \dots \end{aligned}$$

- Neural Network (NN):

$h_{\theta}(\cdot)$  is the output of a NN with weights  $\theta$ .

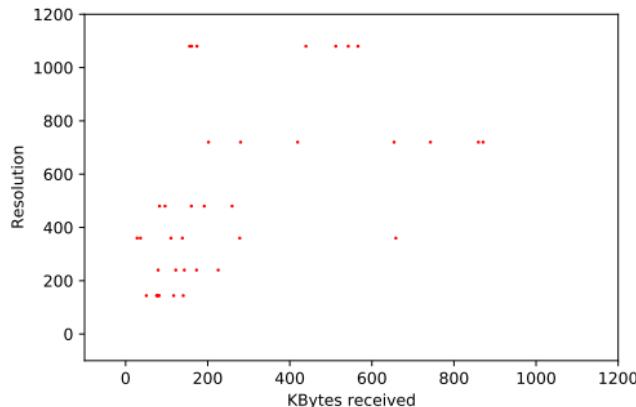
## Training

Find  $\theta^* \triangleq \arg \min_{\theta} J(\theta, \mathbf{X}, \mathbf{y})$

## Example of univariate model

12 / 54

Predict the video resolution just based on KB received in 100ms time slot.  
Are you able to find a simple model  $h(\text{KBytesReceived})$  “by hand”?



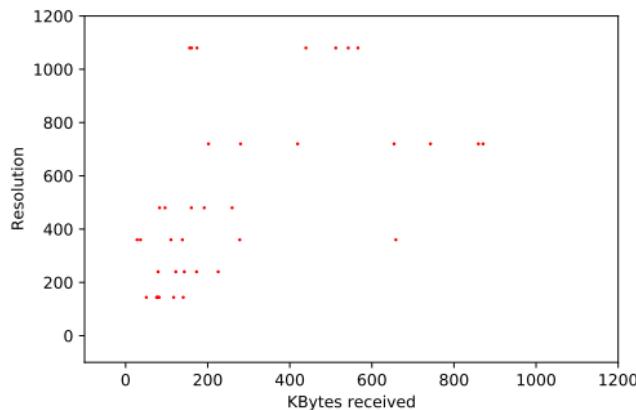
KBytesReceived	avg_qual
10.086	0
16.218	0
16.187	0
13.187	0
36.199	360
16.832	0
17.040	0
16.522	0
77.605	144
16.475	0
16.344	0
75.432	144
11.361	0
11.507	0
17.902	0
278.065	360
19.808	0
27.840	360

[Source of data [[GGA<sup>+</sup>19](#)]]

# Example of univariate model

12 / 54

Predict the video resolution just based on KB received in 100ms time slot.  
Are you able to find a simple model  $h(\text{KBytesReceived})$  “by hand”?



## Generalization

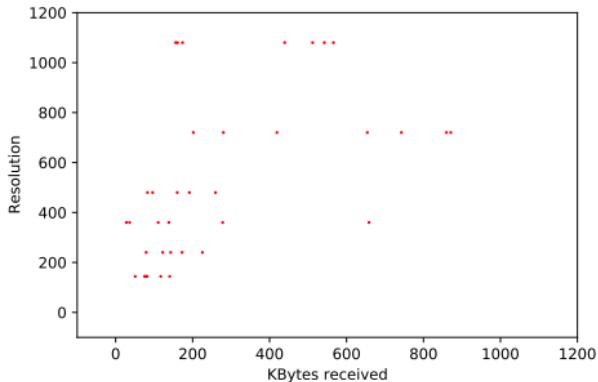
If we see that a new connection with 500 KB / 100ms, what is the predicted video resolution?

[Source of data [GGA<sup>+</sup>19]]

KBytesReceived	avg_qual
10.086	0
16.218	0
16.187	0
13.187	0
36.199	360
16.832	0
17.040	0
16.522	0
77.605	144
16.475	0
16.344	0
75.432	144
11.361	0
11.507	0
17.902	0
278.065	360
19.808	0
27.840	360

# Training a Linear Model: Example

13 / 54



Find  $\theta_0, \theta_1$  such that  $h_{\theta}(\mathbf{x}^{(i)}) \simeq \theta_0 + \theta_1 \cdot x_1^{(i)}$ .

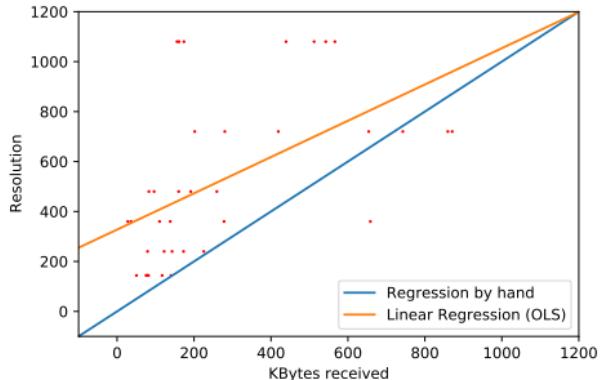
- Regression “by hand”

KBytesReceived	avg_qual
10.086	0
16.218	0
16.187	0
13.187	0
36.199	360
16.832	0
17.040	0
16.522	0
77.605	144
16.475	0
16.344	0
75.432	144
11.361	0
11.507	0
17.902	0
278.065	360
19.808	0
27.840	360

Source of data [GGA<sup>+</sup>19]

# Training a Linear Model: Example

13 / 54



Find  $\theta_0, \theta_1$  such that  $h_{\theta}(\mathbf{x}^{(i)}) \simeq \theta_0 + \theta_1 \cdot x_1^{(i)}$ .

- Regression “by hand”
- Ordinary Least Square (OLS) regression

$$\theta^* = \arg \min_{\theta} J \quad \theta_0^* = 351.8 \quad \theta_1^* = 0.7.$$

KBytesReceived	avg_qual
10.086	0
16.218	0
16.187	0
13.187	0
36.199	360
16.832	0
17.040	0
16.522	0
77.605	144
16.475	0
16.344	0
75.432	144
11.361	0
11.507	0
17.902	0
278.065	360
19.808	0
27.840	360

Source of data [GGA<sup>+</sup>19]

## Bi-variate linear model

14 / 54

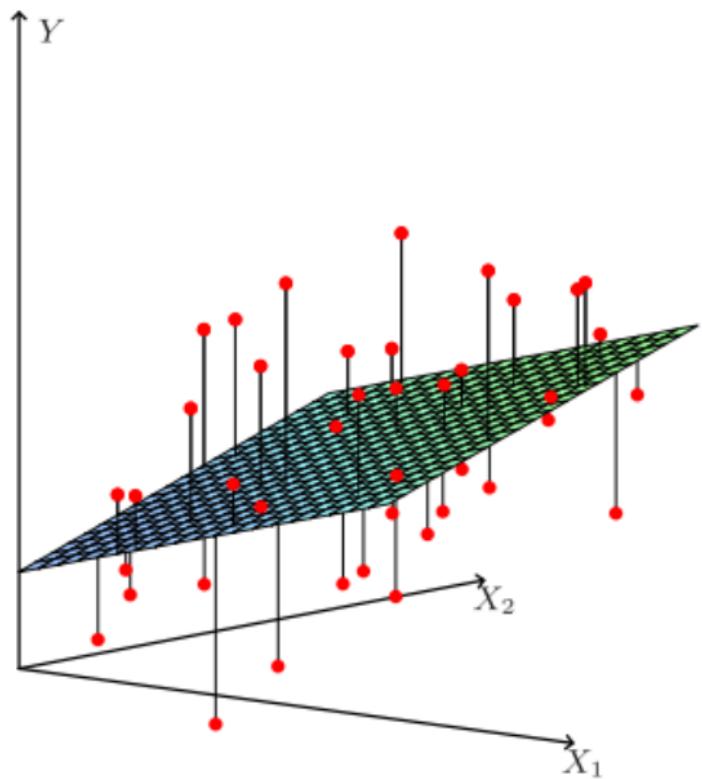


Figure 3.1 of [HTF09]

- Matrix of samples:  $\mathbf{X} = \begin{bmatrix} \mathbf{x}^{(1)T} \\ \vdots \\ \mathbf{x}^{(M)T} \end{bmatrix} = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_N^{(1)} \\ \dots & \dots & \dots & \dots & \dots \\ 1 & x_1^{(M)} & x_2^{(M)} & \dots & x_N^{(M)} \end{bmatrix}$
- Vector of true labels:  $\mathbf{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(M)} \end{bmatrix}$
- Vector of model parameters:  $\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_N \end{bmatrix}$
- Vector of predicted labels:  $\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}^{(1)} \\ \vdots \\ \hat{y}^{(M)} \end{bmatrix} = \begin{bmatrix} h_{\boldsymbol{\theta}}(\mathbf{x}^{(1)}) \\ \vdots \\ h_{\boldsymbol{\theta}}(\mathbf{x}^{(M)}) \end{bmatrix}$

# Examples of loss function

16 / 54

For regression

- Mean Square Error (MSE)

$$J(\boldsymbol{\theta}, \mathbf{X}, \mathbf{y}) = \frac{1}{M} \sum_{i=1}^M \left( \underbrace{y^{(i)} - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})}_{\text{Residual } \varepsilon^{(i)}} \right)^2$$

- Root Mean Square Error (RMSE)

$$J(\boldsymbol{\theta}, \mathbf{X}, \mathbf{y}) = \sqrt{\frac{1}{M} \sum_{i=1}^M (y^{(i)} - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}))^2}$$

- ...

For classification:

- Misclassification Rate

$$J(\boldsymbol{\theta}, \mathbf{X}, \mathbf{y}) = \frac{1}{M} \sum_{i=1}^M I_{y^{(i)} \neq h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})}$$

- ...

## Training

Find  $\theta^* \triangleq \arg \min_{\theta} J(\theta, \mathbf{X}, \mathbf{y})$

Solution algorithm:

- Linear regression: matrix inversion
- Neural network: backpropagation
- ...

Note that

- $\theta^*$  depends on the observed data  $(\mathbf{X}, \mathbf{y})$ .
- If we observed other data  $(\mathbf{X}', \mathbf{y}')$  we would get another  $\theta^*$ .
- ...

## Section 2

### Ordinary Least Squares

For any sample  $\mathbf{x}^{(i)}$ , a linear model is:

$$h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) = \theta_0 \cdot 1 + \theta_1 \cdot x_1^{(i)} + \cdots + \theta_N \cdot x_N^{(i)} = \mathbf{x}^{(i)T} \cdot \boldsymbol{\theta}$$

Assume the loss function below:

$$\text{MSE} = J(\boldsymbol{\theta}, \mathbf{X}, \hat{\mathbf{y}}) = \frac{1}{m} \sum_{i=1}^m \left( y^{(i)} - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \right)^2$$

The model  $h_{\boldsymbol{\theta}^*}$  that minimizes MSE is called **Ordinary Least Squares** (OLS) model.

## Theorem: Normal equation

The optimal parameter vector is

$$\boldsymbol{\theta}^* = (\mathbf{X}^T \cdot \mathbf{X})^{-1} \cdot \mathbf{X}^T \cdot \mathbf{y} \quad (1)$$

(provided that  $(\mathbf{X}^T \cdot \mathbf{X})^{-1}$  is invertible)

NB: Given a certain training set  $(\mathbf{X}, \mathbf{y})$ , we can immediately find the optimal  $\boldsymbol{\theta}^*$ .

# Training a Linear Regression Model (I)

20 / 54

For any sample  $\mathbf{x}^{(i)}$ , the prediction is (This proof is similar to §3.2 of [HTF09].):

$$h_{\theta}(\mathbf{x}^{(i)}) = \theta_0 \cdot 1 + \theta_1 \cdot x_1^{(i)} + \cdots + \theta_N \cdot x_N^{(i)} = \mathbf{x}^{(i)T} \cdot \boldsymbol{\theta}$$

The loss function is:

$$\begin{aligned} J(\boldsymbol{\theta}, \mathbf{X}, \hat{\mathbf{y}}) &= \frac{1}{m} \sum_{i=1}^M \left( y^{(i)} - h_{\theta}(\mathbf{x}^{(i)}) \right)^2 = \frac{1}{m} \sum_{i=1}^M \left( y^{(i)} - \mathbf{x}^{(i)T} \cdot \boldsymbol{\theta} \right)^2 \\ &= \frac{1}{m} (\mathbf{y} - \mathbf{X} \cdot \boldsymbol{\theta})^T \cdot (\mathbf{y} - \mathbf{X} \cdot \boldsymbol{\theta}) \end{aligned}$$

To minimize the function, we set the gradient to 0:

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}, \mathbf{X}, \hat{\mathbf{y}}) = \mathbf{0}$$

Chain rule of derivation:

$$\frac{2}{M} (\mathbf{y} - \mathbf{X} \cdot \boldsymbol{\theta}) \cdot \nabla_{\boldsymbol{\theta}} (\mathbf{y} - \mathbf{X} \cdot \boldsymbol{\theta}) = \mathbf{0} \quad -\mathbf{y} \cdot \mathbf{X} + \mathbf{X} \cdot \boldsymbol{\theta} \cdot \mathbf{X} = \mathbf{0} \quad (4)$$

(2)

$$\frac{2}{M} (\mathbf{y} - \mathbf{X} \cdot \boldsymbol{\theta}) \cdot (-\mathbf{X}) = \mathbf{0} \quad (3)$$

## Training a Linear Regression Model (II)

Let us isolate  $\theta$ :

$$\mathbf{X} \cdot \theta \cdot \mathbf{X} = \mathbf{y} \cdot \mathbf{X} \quad (5)$$

$$\mathbf{X}^T \cdot \mathbf{X} \cdot \theta \cdot \mathbf{X} = \mathbf{X}^T \cdot \mathbf{y} \cdot \mathbf{X} \quad (6)$$

Multiply on the left by  $(\mathbf{X}^T \cdot \mathbf{X})^{-1}$

(Assuming it exists):

$$\theta \cdot \mathbf{X} = (\mathbf{X}^T \cdot \mathbf{X})^{-1} \cdot \mathbf{X}^T \cdot \mathbf{y} \cdot \mathbf{X} \quad (7)$$

Multiply on the right by  $\mathbf{X}^T$

$$\theta \cdot \mathbf{X} \cdot \mathbf{X}^T = (\mathbf{X}^T \cdot \mathbf{X})^{-1} \cdot \mathbf{X}^T \cdot \mathbf{y} \cdot \mathbf{X} \cdot \mathbf{X}^T \quad (8)$$

Multiply on the right by<sup>a</sup>  $(\mathbf{X} \cdot \mathbf{X}^T)^{-1}$ :

$$\theta = (\mathbf{X}^T \cdot \mathbf{X})^{-1} \cdot \mathbf{X}^T \cdot \mathbf{y} \quad (9)$$

---

<sup>a</sup>By the property of transpose matrices, if  $(\mathbf{X}^T \cdot \mathbf{X})^{-1}$  exists, then  $(\mathbf{X} \cdot \mathbf{X}^T)^{-1}$  exists as well.

## Interpreting regression results: sign

22 / 54

You want to predict Resolution based on KBytesReceived:

$$\hat{y}^{(i)} = \theta_0 + \theta_{\text{KBytesReceived}} \cdot x_{\text{KBytesReceived}}^{(i)}$$

KBytesReceived	Resolution
36.199	360
77.605	144
75.432	144
278.065	360

Running the OLS model we get:

	coef	P> t
const	351.7606	0.000
KBytesReceived	0.7115	0.001

The model is

$$\hat{y}^{(i)} = 351.76 + 0.71 \cdot x_{\text{KBytesReceived}}^{(i)}$$

Does the **sign** makes sense?  
Positive/negative dependency.

*p*-values: *significance* of a coefficient.

	coef	P> t
const	351.7606	0.000
KBytesReceived	0.7115	0.001

Hypothesis testing:

- **Null hypothesis:** No dependency between KBytesReceived and the target.

*p*-values: *significance* of a coefficient.

	coef	P> t
const	351.7606	0.000
KBytesReceived	0.7115	0.001

Hypothesis testing:

- **Null hypothesis:** No dependency between KBytesReceived and the target.

- Performing OLS, under the null hp

$$\mathbb{P}\left[|\hat{\theta}_{\text{KBytesReceived}}^*| \geq 0.71\right] = 0.001$$

- Would you *accept* or *reject* the null hp?

*p*-values: *significance* of a coefficient.

	coef	P> t
const	351.7606	0.000
KBytesReceived	0.7115	0.001

Hypothesis testing:

- **Null hypothesis:** No dependency between KBytesReceived and the target.

- Performing OLS, under the null hp

$$\mathbb{P}\left[|\hat{\theta}_{\text{KBytesReceived}}^*| \geq 0.71\right] = 0.001$$

- Would you *accept* or *reject* the null hp?

- Accept  $\Rightarrow 0.71$  is not significant
  - Reject  $\Rightarrow 0.71$  is significant.

$p$ -values: *significance* of a coefficient.

	coef	P> t
const	351.7606	0.000
KBytesReceived	0.7115	0.001

Hypothesis testing:

- **Null hypothesis:** No dependency between KBytesReceived and the target.

- Performing OLS, under the null hp

$$\mathbb{P}\left[|\hat{\theta}_{\text{KBytesReceived}}^*| \geq 0.71\right] = 0.001$$

- Would you *accept* or *reject* the null hp?

- Accept  $\Rightarrow 0.71$  is not significant
- Reject  $\Rightarrow 0.71$  is significant.

Usually:

- $p$ -value  $\leq 5\% \Rightarrow$  reject null hp  $\Rightarrow$  Coeff significant
- o.w. Coeff not significant

# Interpreting regression results

24 / 54

Predict Resolution based on  
packetsSent, KBytesReceived and  
BufferHealth:

$$\hat{y}^{(i)} = \theta_0 + \theta_1 \cdot x_1^{(i)} + \theta_2 \cdot x_2^{(i)} + \theta_3 \cdot x_3^{(i)}$$

packetsSent	KBytesReceived	BufferHealth	Resolution
6	36.199	10.241165	360
10	77.605	4.446780	144
14	75.432	3.989780	144
33	278.065	3.700462	360
6	27.840	4.512780	360
58	658.375	9.454706	360
14	77.429	4.606780	144
33	201.903	5.301853	720
18	172.740	3.638107	240
23	181.476	5.314732	240

Training the OLS model:

	coef	P> t
const	282.8794	0.017
packetsSent	-0.5551	0.925
KBytesReceived	0.5986	0.194
BufferHealth	18.3779	0.343

The model is

$$\hat{y}^{(i)} = 282.9 - 0.56 \cdot x_1^{(i)} + 0.60 \cdot x_2^{(i)} + 18.34 \cdot x_3^{(i)}$$

Interpretation:

- Fixing all features
- 1KB more received (in the 100 ms window)  $\Rightarrow$  resolution  $\nearrow$  by 0.60p
- $\Leftrightarrow$  increase of 1MB  $\Rightarrow$  resolution  $\nearrow$  of 600pR.

# Interpreting regression results

25 / 54

		coef	P> t
	<b>const</b>	282.8794	0.017
	<b>PacketsSent</b>	-0.5551	0.925
	<b>KBytesReceived</b>	0.5986	0.194
	<b>BufferHealth</b>	18.3779	0.343

$$\hat{y}^{(i)} = 282.9 - 0.56 \cdot x_1^{(i)} + 0.60 \cdot x_2^{(i)} + 18.34 \cdot x_3^{(i)}$$

Do the **signs** make sense?

		coef	P> t
	<b>const</b>	282.8794	0.017
	<b>PacketsSent</b>	-0.5551	0.925
	<b>KBytesReceived</b>	0.5986	0.194
	<b>BufferHealth</b>	18.3779	0.343

$$\hat{y}^{(i)} = 282.9 - 0.56 \cdot x_1^{(i)} + 0.60 \cdot x_2^{(i)} + 18.34 \cdot x_3^{(i)}$$

Do the **signs** make sense?

Are coefficients **significant**?

## Theorem

Assume that

- for any sample  $\mathbf{x}$ , the target is a random variable (r.v.):

$$\hat{y} = \mathbf{x}^T \cdot \boldsymbol{\beta} + \epsilon \quad (10)$$

- $\boldsymbol{\beta}$ : some “true” parameter vector
- $\epsilon$ : noise, Gaussian r.v.  $\epsilon \sim N(0, \sigma^2)$
- $\sigma^2$ : variance of the residuals

- The true label  $y^{(i)}$  is a realization of  $\hat{y}$ .

Then

- the coefficients  $\hat{\theta}_j^*$  we get from linear regression are also Gaussian r.v. They are **unbiased**, i.e.,

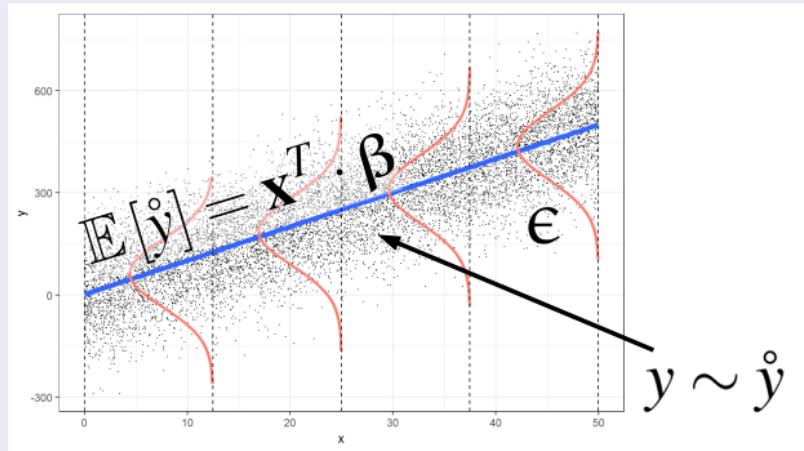
$$\mathbb{E} \left[ \hat{\theta}_j^* \right] = \beta_j$$

# How is the $p$ -value computed (II)

28 / 54

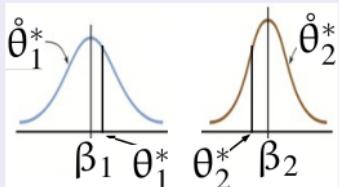
## Theorem

Assume that



Picture from [LR19], Ch. I

Then



## How is the $p$ -value computed (III)

29 / 54

Proof of the previous claim.<sup>1</sup>

- Let's write (10) in vectorial form. Given a dataset  $\mathbf{X} = \begin{bmatrix} \mathbf{x}^{(1)\top} \\ \vdots \\ \mathbf{x}^{(M)\top} \end{bmatrix}$ , the target vector is:

$$\hat{\mathbf{y}} = \mathbf{X} \cdot \boldsymbol{\beta} + \boldsymbol{\epsilon}$$

It is a Gaussian r.v., since it is a constant matrix  $\mathbf{X} \cdot \boldsymbol{\beta}$  plus a random vector  $\boldsymbol{\epsilon}$ .

- The parameter vector we get from OLS regression is:

$$\hat{\boldsymbol{\theta}}^* = (\mathbf{X}^T \cdot \mathbf{X})^{-1} \cdot \mathbf{X}^T \cdot \hat{\mathbf{y}}$$

It is a Gauss.r.v., since it is a constant matrix  $(\mathbf{X}^T \cdot \mathbf{X})^{-1} \cdot \mathbf{X}^T$  multiplied by the Gauss.r.v.  $\hat{\mathbf{y}}$ .

- The mean is

$$\mathbb{E} [\hat{\boldsymbol{\theta}}^*] = (\mathbf{X}^T \cdot \mathbf{X})^{-1} \cdot \mathbf{X}^T \cdot \mathbb{E} [\hat{\mathbf{y}}] = (\mathbf{X}^T \cdot \mathbf{X})^{-1} \cdot \mathbf{X}^T \cdot \mathbf{X} \cdot \boldsymbol{\beta} = \boldsymbol{\beta}$$

and thus  $\mathbb{E} [\hat{\theta}_j^*] = \beta_j$

- We can also compute the variance of  $\hat{\boldsymbol{\theta}}^*$ .<sup>2</sup>

<sup>1</sup>To know more, check Sec. 3.8 of [HTF09], [these videos](#) and [this video](#).

<sup>2</sup>

## How is the $p$ -value computed (IV)

- Compute

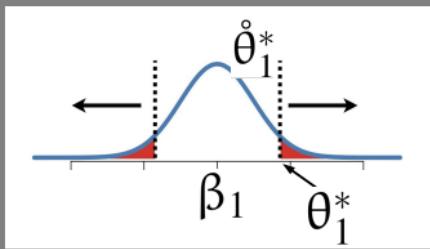
$$\boldsymbol{\theta}^* = (\mathbf{X}^T \cdot \mathbf{X})^{-1} \cdot \mathbf{X}^T \cdot \mathbf{y}$$

- $\boldsymbol{\theta}^*$  is a realization of the Gauss.r.v.

$$\hat{\boldsymbol{\theta}}^* = (\mathbf{X}^T \cdot \mathbf{X})^{-1} \cdot \mathbf{X}^T \cdot \hat{\mathbf{y}}$$

- For any feature  $j$ :

- $\theta_j^*$  realization of  $\hat{\theta}_j^*$ , such that  $\mathbb{E} [\hat{\theta}_j^*] = \beta_j$
- Assume true parameter  $\beta_j = 0$ .
- Compute the variance  $\text{Var}(\hat{\theta}_j^*)$ .
- Draw the probability density of  $\hat{\theta}_j^*$ .



p-value
$\mathbb{P} [ \hat{\theta}_j^*  \geq \theta_j^*]$

# Significance degradation

Adding features can degrade significance

	coef	P> t
<b>const</b>	351.7606	0.000
<b>KBytesReceived</b>	0.7115	0.001

	coef	P> t
<b>const</b>	282.8794	0.017
<b>PacketsSent</b>	-0.5551	0.925
<b>KBytesReceived</b>	0.5986	0.194
<b>BufferHealth</b>	18.3779	0.343

## Theorem: prediction and true value averages

In an OLS regression, the average of the predictions equals the average of the true values of the target:

$$\bar{\hat{y}} = \bar{y}$$

where:

$$\bar{\hat{y}} \triangleq \frac{1}{M} \sum_{i=1}^M \hat{y}^{(i)} \quad \bar{y} \triangleq \frac{1}{M} \sum_{i=1}^M y^{(i)}$$

Proof

Normal equation:

$$\theta^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

## Theorem: prediction and true value averages

In an OLS regression, the average of the predictions equals the average of the true values of the target:

$$\bar{\hat{y}} = \bar{y}$$

where:

$$\bar{\hat{y}} \triangleq \frac{1}{M} \sum_{i=1}^M \hat{y}^{(i)} \quad \bar{y} \triangleq \frac{1}{M} \sum_{i=1}^M y^{(i)}$$

Proof

Normal equation:

$$\theta^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Multiplying on the left by  $\mathbf{X}^T \mathbf{X}$ :

$$\mathbf{X}^T \mathbf{X} \cdot \underbrace{\theta^*}_{\mathbf{y}} = \mathbf{X}^T \mathbf{y}$$

## Theorem: prediction and true value averages

In an OLS regression, the average of the predictions equals the average of the true values of the target:

$$\bar{\hat{y}} = \bar{y}$$

where:

$$\bar{\hat{y}} \triangleq \frac{1}{M} \sum_{i=1}^M \hat{y}^{(i)} \quad \bar{y} \triangleq \frac{1}{M} \sum_{i=1}^M y^{(i)}$$

Proof

Normal equation:

$$\theta^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Multiplying on the left by  $\mathbf{X}^T \mathbf{X}$ :

$$\mathbf{X}^T \underbrace{\mathbf{X} \cdot \theta^*}_{\hat{y}} = \mathbf{X}^T \mathbf{y}$$

## Theorem: prediction and true value averages

In an OLS regression, the average of the predictions equals the average of the true values of the target:

$$\bar{\hat{y}} = \bar{y}$$

where:

$$\bar{\hat{y}} \triangleq \frac{1}{M} \sum_{i=1}^M \hat{y}^{(i)} \quad \bar{y} \triangleq \frac{1}{M} \sum_{i=1}^M y^{(i)}$$

Proof

Normal equation:

$$\theta^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Multiplying on the left by  $\mathbf{X}^T \mathbf{X}$ :

$$\mathbf{X}^T \mathbf{X} \cdot \underbrace{\theta^*}_{\hat{y}} = \mathbf{X}^T \mathbf{y} \implies \mathbf{X}^T \hat{y} = \mathbf{X}^T \mathbf{y}$$

## Theorem: prediction and true value averages

In an OLS regression, the average of the predictions equals the average of the true values of the target:

$$\bar{\hat{y}} = \bar{y}$$

where:

$$\bar{\hat{y}} \triangleq \frac{1}{M} \sum_{i=1}^M \hat{y}^{(i)} \quad \bar{y} \triangleq \frac{1}{M} \sum_{i=1}^M y^{(i)}$$

Proof

Normal equation:

$$\theta^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

The first element is:

$$[1 \dots 1] \cdot \hat{\mathbf{y}} = [1 \dots 1] \cdot \mathbf{y}$$

Multiplying on the left by  $\mathbf{X}^T \mathbf{X}$ :

$$\mathbf{X}^T \underbrace{\mathbf{X} \cdot \theta^*}_{\hat{\mathbf{y}}} = \mathbf{X}^T \mathbf{y} \implies \mathbf{X}^T \hat{\mathbf{y}} = \mathbf{X}^T \mathbf{y}$$

$$\hat{y} = \theta_0 + \underbrace{\theta_1 x_1 + \cdots + \theta_N x_N}_{\text{dependency}}$$

- $\theta_0$  does not capture  $(x, y)$  dependency
- $\theta_0$  just “aligns” predictions to meet the label average.

## Section 3

### **Validation of a model**

- Supervised learning: we construct a model  $h(\mathbf{x})$  based on observed  $(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots)$  for the purpose of approximating  $y$  for new samples  $\mathbf{x}$ .
- **Generalization:**<sup>a</sup> when the constructed model  $h(\mathbf{x})$  is good at approximating labels  $y$  of **new** samples  $\mathbf{x}$ :

$$h(\mathbf{x}) \simeq y$$

- How do we test if the model *generalizes*.



- Student analogy
  - He/she checks the answers during **training**
  - Answers are hidden during the **test**

---

<sup>a</sup>Sec. 3.1 of [GBD92].

# Training and test sets

36 / 54

BufferHealth	BufferProgress	BufferValid	label	label_num
10.241165	0.015357	true	q360p	360
4.446780	0.007103	true	q144p	144
3.989780	0.006509	true	q144p	144
3.700462	0.005897	true	q360p	360
4.512780	0.007156	true	q360p	360
9.454706	0.016805	true	q360p	360
4.606780	0.008046	true	q144p	144
5.301853	0.007990	true	q720p	720
3.638107	0.005493	true	q240p	240
5.314732	0.007407	true	q240p	240
8.554780	0.014515	true	q480p	480
4.189780	0.007516	true	q360p	360
3.633641	0.005897	true	q480p	480
1.495841	0.002473	true	q720p	720
8.802211	0.014076	true	q1080p	1080
4.611142	0.009263	true	q144p	144
5.590378	0.009113	true	q480p	480
4.940168	0.008851	true	q1080p	1080
4.940168	0.008851	true	q1080p	1080
9.239532	0.016335	true	q720p	720

X y

# Training and test sets

36 / 54

BufferHealth	BufferProgress	BufferValid	label	label_num
10.241165	0.015357	true	q360p	360
4.446780	0.007103	true	q144p	144
3.989780	0.006509	true	q144p	144
3.700462	0.005897	true	q360p	360
4.512780	0.007156	true	q360p	360
9.454780	0.016805	true	q360p	360
4.606780	0.008046	true	q144p	144
5.301845	0.007990	true	q720p	720
3.638107	0.005193	true	q360p	360
5.314732	0.009400	true	q240p	240
8.554780	0.011688	true	q480p	480
4.189780	0.007516	true	q360p	360
3.633641	0.005897	true	q480p	480
1.495841	0.002473	true	q720p	720
8.802211	0.014076	true	q1080p	1080
4.611142	0.009263	true	q144p	144
5.590378	0.009113	true	q480p	480
4.940168	0.008851	true	q1080p	1080
4.940168	0.008551	true	q1080p	1080
9.239532	0.016553	true	q720p	720

X train

Y train

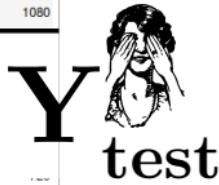
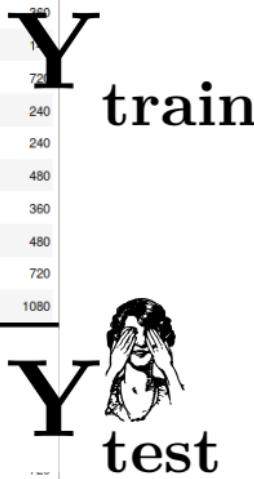
X test

Y test

# Training and test sets

36 / 54

BufferHealth	BufferProgress	BufferValid	label	label_num
10.241165	0.015357	true	q360p	360
4.446780	0.007103	true	q144p	144
3.989780	0.006509	true	q144p	144
3.700462	0.005897	true	q360p	360
4.512780	0.007156	true	q360p	360
9.454770	0.016805	true	q360p	360
4.606780	0.008046	true	q144p	144
5.301950	0.07790	true	q720p	720
3.638107	0.00493	true	q240p	240
5.314732	0.009400	true	q240p	240
8.554780	0.011688	true	q480p	480
4.189780	0.007516	true	q360p	360
3.633641	0.005897	true	q480p	480
1.495841	0.002473	true	q720p	720
8.802211	0.014076	true	q1080p	1080
4.611142	0.009263	true	q144p	144
5.590378	0.009113	true	q480p	480
4.940168	0.008851	true	q1080p	1080
4.940168	0.008551	true	q1080p	1080
9.239532	0.01635	true	q720p	720



# Training and test sets

36 / 54

Train using only the **training set**  
 $(\mathbf{X}, \hat{\mathbf{y}})$ :

$$\boldsymbol{\theta}^* \triangleq \arg \min_{\boldsymbol{\theta}} J(\boldsymbol{\theta}, \mathbf{X}_{\text{train}}, \mathbf{y}_{\text{train}})$$

The trained model is

$$h_{\boldsymbol{\theta}^*}(\cdot)$$

Evaluate the quality of the trained model via:

$$J(\boldsymbol{\theta}^*, \mathbf{X}_{\text{test}}, \mathbf{y}_{\text{test}})$$

*(Test error or generalization error<sup>a</sup> or out-of-sample error)*

The  $J$  used during training and test may be different.

BufferHealth	BufferProgress	BufferValid	label	label_num
10.241165	0.015357	true	q360p	360
4.446780	0.007103	true	q144p	144
3.989780	0.006509	true	q144p	144
3.700462	0.005897	true	q360p	360
4.512780	0.007156	true	q360p	360
9.454770	0.016805	true	q360p	360
4.606780	0.008046	true	q144p	144
5.301950	0.07790	true	q720p	720
3.638107	0.00493	true	q240p	240
5.314732	0.009400	true	q240p	240
8.554780	0.011688	true	q480p	480
4.189780	0.007516	true	q360p	360
3.633641	0.005897	true	q480p	480
1.495841	0.002473	true	q720p	720
8.802211	0.014076	true	q1080p	1080
4.611142	0.009263	true	q144p	144
5.590378	0.009113	true	q480p	480
4.940168	0.008851	true	q1080p	1080
4.940168	0.008551	true	q1080p	1080
9.239532	0.01635	true	q720p	720

**X** train      **Y** train

**X** test      **Y** test



Which method is better?

BufferHealth	BufferProgress	BufferValid	label	label_num
10.241165	0.015357	true	q360p	360
4.446780	0.007103	true	q144p	144
3.989780	0.006509	true	q144p	144
3.700462	0.005897	true	q360p	360
4.512780	0.007156	true	q360p	360
9.454706	0.016805	true	q360p	360
4.606780	0.008046	true	q144p	144
5.301853	0.007990	true	q720p	720
3.638107	0.005493	true	q240p	240
5.314732	0.009400	true	q240p	240
8.554780	0.011688	true	q480p	480
4.189780	0.007516	true	q360p	360
3.633641	0.005897	true	q480p	480
1.495641	0.002473	true	q720p	720
8.802211	0.014076	true	q1080p	1080
4.611142	0.009263	true	q144p	144
5.590378	0.009113	true	q480p	480
4.940168	0.008851	true	q1080p	1080
4.940168	0.008851	true	q1080p	1080
9.239532	0.016335	true	q720p	720

BufferHealth	BufferProgress	BufferValid	label	label_num
10.241165	0.015357	true	q360p	360
4.446780	0.007103	true	q144p	144
3.989780	0.006509	true	q144p	144
3.700462	0.005897	true	q360p	360
4.512780	0.007156	true	q360p	360
9.454706	0.016805	true	q360p	360
4.606780	0.008046	true	q144p	144
5.301853	0.007990	true	q720p	720
3.638107	0.005493	true	q240p	240
5.314732	0.009400	true	q240p	240
8.554780	0.011688	true	q480p	480
4.189780	0.007516	true	q360p	360
3.633641	0.005897	true	q480p	480
1.495641	0.002473	true	q720p	720
8.802211	0.014076	true	q1080p	1080
4.611142	0.009263	true	q144p	144
5.590378	0.009113	true	q480p	480
4.940168	0.008851	true	q1080p	1080
4.940168	0.008851	true	q1080p	1080
9.239532	0.016335	true	q720p	720

What if the data provider first experimented with all low resolutions and finally with all the highest?



Go to notebook 02.regression/a.regression.ipynb

## Section 4

### **Instability of a model (Variance)**

$$\text{Var}(j) = \frac{1}{M-1} \sum_{i=1}^M (x_j^{(i)} - \mu_j)^2$$

- Suppose  $x_j^{(i)} = c$  for all samples  $i$
- $\mathbf{X} = \begin{bmatrix} 1 & x_1^{(1)} & \dots & c & \dots & x_N^{(1)} \\ \vdots & \vdots & & \vdots & \vdots & \vdots \\ 1 & x_1^{(M)} & \dots & c & \dots & x_N^{(M)} \end{bmatrix}$
- Not full rank
- $\Rightarrow (\mathbf{X}^T \mathbf{X})^{-1}$  does not exist
- $\theta^* = (\mathbf{X}^T \cdot \mathbf{X})^{-1} \cdot \mathbf{X}^T \cdot \mathbf{y}$  impossible

- Two features  $j_1, j_1$  (columns of  $\mathbf{X}$ ) are **collinear** if they are proportional

$$\begin{bmatrix} x_{j_1}^{(1)} \\ \vdots \\ x_{j_1}^{(N)} \end{bmatrix} = \alpha \begin{bmatrix} x_{j_2}^{(1)} \\ \vdots \\ x_{j_2}^{(N)} \end{bmatrix}$$

In this case, column  $j_2$  **adds no information** about prediction.

- $\Rightarrow \mathbf{X}$  has no full rank
- $\Rightarrow (\mathbf{X}^T \cdot \mathbf{X})$  is not invertible
- $\Rightarrow$  not unique.
- $\Rightarrow$  The normal equation

$$\boldsymbol{\theta}^* = (\mathbf{X}^T \cdot \mathbf{X})^{-1} \cdot \mathbf{X}^T \cdot \mathbf{y}$$

cannot be computed.

- In real world, perfect collinearity is rare, but

$$\begin{bmatrix} x_{j_1}^{(1)} \\ \vdots \\ x_{j_1}^{(N)} \end{bmatrix} \simeq \alpha \begin{bmatrix} x_{j_2}^{(1)} \\ \vdots \\ x_{j_2}^{(N)} \end{bmatrix}$$

and the computer is able to compute the normal equation, but ...

- For any model  $h_\theta$ :

$$\begin{aligned} h_\theta(\mathbf{x}) &= \theta_{j_1}x_{j_1} + \theta_{j_2}x_{j_2} + \sum_{j \notin \{j_1, j_2\}} \theta_j x_j \\ &\simeq (\alpha\theta_{j_1} + \theta_{j_2})x_{j_2} + \sum_{j \notin \{j_1, j_2\}} \theta_j x_j \end{aligned}$$

- Infinite pairs of  $(\theta_{j_1}, \theta_{j_2})$  would be almost equivalent
  - Ex.

- How does OLS choose among such pairs? Randomly
- Small differences of  $\mathbf{X} \Rightarrow$  big differences in  $(\theta_{j_1}^*, \theta_{j_2}^*)$
- $\Rightarrow$  Variance

- If  $j$ -th feature is a linear combination of others

$$\mathbf{X}_j = \alpha_1 \mathbf{X}_{j_1} + \alpha_2 \mathbf{X}_{j_2} + \dots$$

- $\mathbf{X} = [\mathbf{1} | \mathbf{X}_1 | \dots | \mathbf{X}_N]$  no full rank
- The normal equation

$$\boldsymbol{\theta}^* = (\mathbf{X}^T \cdot \mathbf{X})^{-1} \cdot \mathbf{X}^T \cdot \mathbf{y}$$

cannot be computed.

- How to discover Multi-collinearity? ...

- Humans are good at generalizing knowledge ...
- ... because their perception models have low variance

## Variance

A model suffers **high variance** if, by perturbing a bit the training dataset, the model changes completely.

Suppose  $\tilde{\mathbf{X}}, \tilde{\mathbf{y}}$  is a slightly perturbed version of  $\mathbf{X}, \mathbf{y}$ . If a model has high variance:

$$\tilde{\theta}^* = \arg \min_{\theta} J(\theta, \tilde{\mathbf{X}}, \tilde{\mathbf{y}})$$

completely different than

$$\theta^* = \arg \min_{\theta} J(\theta, \mathbf{X}, \mathbf{y})$$

- **Variance** is the contrary of **Stability**
- High variance  $\Rightarrow$  high sensitivity to training data.  
(§5 of [BK11])





Go to notebook 02regression/.a.regression.ipynb

## Section 5

### **Cross-validation**

How can we be sure that, if we change train/test split the error does not change?  
⇒ **Cross-validation.**

---

## Algorithm 1 $K$ -fold validation

---

- 1: Divide the dataset in  $K$  subsets
  - 2: **for**  $i = 1$  to  $K$  **do**
  - 3:     Keep subset  $i$  for test
  - 4:     Train on all the others
  - 5:     Compute test error
  - 6: **end for**
  - 7: Error = average of test errors
- 

BufferHealth	BufferProgress	BufferValid	label	label_num
10.241165	0.015357	true	q360p	360
4.446780	0.007103	true	q144p	144
3.989780	0.006509	true	q144p	144
3.700462	0.005897	true	q360p	360
4.512780	0.007156	true	q360p	360
9.454706	0.016805	true	q360p	360
4.606780	0.008046	true	q144p	144
5.301853	0.007990	true	q720p	720
3.638107	0.005493	true	q240p	240
5.314732	0.009400	true	q240p	240
8.554780	0.011688	true	q480p	480
4.189780	0.007516	true	q360p	360
3.633641	0.005897	true	q480p	480
1.495841	0.002473	true	q720p	720
8.802211	0.014076	true	q1080p	1080
4.611142	0.009263	true	q144p	144
5.590378	0.009113	true	q480p	480
4.940168	0.008851	true	q1080p	1080
4.940168	0.008851	true	q1080p	1080
9.239532	0.016335	true	q720p	720

TRAINING SET

TEST SET



Go to notebook 02.regression/a.regression.ipynb

### In this lesson

- Supervised Learning
- First Model in Python  
(Linear Regression)
- Feature selection
- Validation: Train/Test; Cross-validation

## Regression (continued)

- Polynomial Regression
- Model Variance / Complexity
- Regularization
- Scaling
- Feature Selection

## Classification

- Logistic Regression
- Classification Performance
- Class imbalance

- Dealing with Multi-collinearity [DEB<sup>+</sup>13]
- Variance Inflation Factor for testing Multi-collinearity: pag. 101-102 of [JWHT13]
- Cross-validation from [machinelearningmastery](#)

-  Eric Bauer and Ron Kohavi, *An Empirical Comparison of Voting Classification Algorithms : Bagging , Boosting , and Variants*, Machine Learning **38** (2011), no. 1998, 1–38.
-  Carsten F. Dormann, Jane Elith, Sven Bacher, Carsten Buchmann, Gudrun Carl, Gabriel Carré, Jaime R. García Marquéz, Bernd Gruber, Bruno Lafourcade, Pedro J. Leitão, Tamara Münkemüller, Colin McClean, Patrick E. Osborne, Björn Reineking, Boris Schröder, Andrew K. Skidmore, Damaris Zurell, and Sven Lautenbach, *Collinearity: A review of methods to deal with it and a simulation study evaluating their performance*, Ecography **36** (2013), no. 1, 027–046.
-  Stuart Geman, Elie Bienenstock, and René Doursat, *Neural Networks and the Bias/Variance Dilemma*, Neural Computation **4** (1992), no. 1, 1–58.

-  Craig Gutterman, Katherine Guo, Sarthak Arora, Xiaoyang Wang, Les Wu, Ethan Katz-Bassett, and Gil Zussman, *Request: Real-time QoE detection for encrypted YouTube traffic*, ACM MMSys, 2019.
-  Trevor Hastie, Robert Tibshirani, and Jerome Friedman, *The Elements of Statistical Learning*, 2 ed., vol. 1, Springer, 2009.
-  Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani, *An introduction to Statistical Learning*, vol. 7, 2013.
-  Muhammad Jawad Khokhar, Thibaut Ehlinger, and Chadi Barakat, *From Network Traffic Measurements to QoE for Internet Video*, IFIP Networking, 2019, pp. 1–9.
-  Julie Legler and Paul Roback, *Broadening Your Statistical Horizons*, online, 2019.

-  Simon Thorpe, *Intelligence artificielle et intelligence naturelle – vers l'IA bio-inspiré*, <https://www.college-de-france.fr/site/stephane-mallat/seminar-2020-02-26-11h00.htm>,  
2020.