

#### Aprendizaje automático: riesgos computacionales y sociales

Universidad de Oviedo Academia Asturiana de Ciencia e Ingeniería



Antonio Bahamonde



Universidad de Oviedo

#### Introducción

- Aprendizaje automático: un punto de vista actual
- Riesgos
  - Computacionales: generalización, sesgo (bias)-varianza, regularización
  - Sociales: laborales, empresariales, humanos





#### Contents

- Linear regression
  - Gradient descent
  - Probabilistic interpertation
- Generalization
- Bias and Variance
- Regularization
- Riesgos sociales del AA (LLM)
- Reference (main)
  - Main notes, CS229, Stanford University (Andrew Ng)

Living area (feet <sup>2</sup> )	#bedrooms	Price $(1000$ \$s)
2104	3	400
1600	3	330
2400	3	369
1416	2	232
3000	4	540
• • •	• • •	• • •

 $\theta_i$ 's are the parameters (also called weights) parameterizing the space of linear functions mappings

 $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$ 

 $h: \mathcal{X} \mapsto \mathcal{Y}$ 

When there is no risk of confusion, we will drop the  $\theta$  subscript in  $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$ 

write it more simply as h(x). To simplify our notation, we also introduce the convention of letting  $x_0 = 1$  (this is the **intercept** term), so that

$$h(x) = \sum_{i=0}^{a} \theta_i x_i = \theta^T x,$$

7

Given a training set, how do we pick, or learn, the parameters  $\theta$ ? One reasonable method seems to be to make h(x) close to y, at least for the training examples we have.

close the h(x)'s are to the corresponding y's. We define the cost function:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^{n} (h_{\theta}(x^{(i)}) - y^{(i)})^2.$$

To formalize this, we will define a function that measures, for each value of the  $\theta$ 's, how

least-squares cost function that gives rise to the ordinary least squares regression model



 $\mathbf{Y}$ 

-3-1.0 -0.50.0 0.51.0 $x_1$ 

#### Linear Regression Optimization of $\theta$ 0.550.500.45 $MSE^{(train)}$ 0.40 0.350.30 **J(θ)** 0.250.20

0.51.51.0 $w_1$ 

 $y \approx \theta_1 x_1$ 

#### Gradient descent

#### Gradient descent

Given a function  $J(\theta 1, \theta 2)$  we look for the values of  $\theta 1$  and  $\theta 2$ that make the value of J to a minimum

Algorithm:

- Assign arbitrary values to  $\theta$ 1 and  $\theta$ 2
- 0

Modify the values of 91 and 92 to reduce the value of J Until we believe that we have reached a minimum



#### Gradient descent

To solve

GD is the algorithm

 $\theta \leftarrow \theta_0$ 

repeat

return  $\theta$ 





#### Linear r



- $\frac{\partial}{\partial \theta_j} J(\theta) = \frac{\partial}{\partial \theta_j} \frac{1}{2} \left( h_{\theta} \right)$  $= 2 \cdot \frac{1}{2} \left( h_{\theta} \right)$ 

  - $= (h_{\theta}(x) -$
  - $= (h_{\theta}(x) -$

$$\begin{aligned} & \theta_{j} := \theta_{j} - \alpha \frac{\partial}{\partial \theta_{j}} J(\theta). \\ & \frac{\partial}{\partial \theta_{j}} \frac{1}{2} (h_{\theta}(x) - y)^{2} \\ & 2 \cdot \frac{1}{2} (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_{j}} (h_{\theta}(x) - y) \\ & (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_{j}} \left( \sum_{i=0}^{d} \theta_{i} x_{i} - y \right) \\ & (h_{\theta}(x) - y) x_{j} \end{aligned}$$

 $\theta_j := \theta_j + \alpha \left( y^{(i)} - h_\theta(x^{(i)}) \right) x_j^{(i)}.$ 



$$\theta_j := \theta_j + \alpha \left( y^{(i)} - h_\theta(x^{(i)}) \right) x_j^{(i)}.$$

the magnitude of the update is proportional to the error term  $(y^{(i)} - h_{\theta}(x^{(i)}))$ :

- if we are encountering a training example on which our prediction nearly matches the actual value of y<sup>(i)</sup>, then we find that there is little need to change the parameters;
- a larger change to the parameters will be made if our prediction has a large error

For a single training example, this gives the LMS ("least mean squares") update rule



#### Probabilistic interpretation

$$y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)},$$

$$p(y^{(i)}|x^{(i)};\theta) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right).$$

$$L(\theta) = \prod_{i=1}^{n} p(y^{(i)} \mid x^{(i)}; \theta)$$
  
= 
$$\prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(y^{(i)} - \theta^{T} x^{(i)})^{2}}{2\sigma^{2}}\right)$$

$$p(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(\epsilon^{(i)})^2}{2\sigma^2}\right)$$

$$(\theta); \theta)$$

#### Probabilistic interpretation

$$\begin{split} \ell(\theta) &= \log L(\theta) \\ &= \log \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^{T} x^{(i)})^{2}}{2\sigma^{2}}\right) \\ &= \sum_{i=1}^{n} \log \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^{T} x^{(i)})^{2}}{2\sigma^{2}}\right) \\ &= n \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{\sigma^{2}} \cdot \frac{1}{2} \sum_{i=1}^{n} (y^{(i)} - \theta^{T} x^{(i)})^{2}. \end{split}$$

#### Probabilistic interpretation

Least-squares regression corresponds to finding the maximum likelihood estimate of  $\theta$ .

 $\theta$  did not depend on what was  $\sigma$ , even if  $\sigma$  were unknown.

$$\frac{1}{2} \sum_{i=1}^{n} (y^{(i)})$$

$$-\theta^T x^{(i)})^2$$

- unseen test example
- to as the training loss/error/cost.

The generalization of machine learning models: their performances on

 Supervised learning problems, given a training dataset, we typically learn a model  $h_{\theta}$  by minimizing a loss/cost function J( $\theta$ ), which encourages  $h_{\theta}$  to fit the data. This loss function for training purposes is oftentimes referred

- Minimizing the training loss is not our ultimate goal—it is merely our important evaluation metric of a model is the loss on unseen test examples, which is oftentimes referred to as the test error.
  - minimizing the training error may not always lead to a small test error 0
  - 0

approach towards the goal of learning a predictive model. The most

the model overfits the data if the model predicts accurately on the training dataset but doesn't generalize well to other test examples, that is, if the training error is small but the test error is large.

Test distribution D. The expected loss/error over the randomness of the test example is called the test loss/error

 $L(\theta) = \mathbb{E}_{(x,y)}$ 

model parameterizations.

We will decompose the test error into "bias" and "variance" and their tradeoffs.

We will discuss when overfitting and underfitting will occur and be avoided.

$$_{)\sim\mathcal{D}}[(y-h_{\theta}(x))^{2}]$$

- We'll see how the test error is influenced by the learning procedure, especially the choice of

- The training inputs x<sup>(i)</sup>'s are randomly chosen and the outputs y<sup>(i)</sup> are generated by y<sup>(i)</sup>= h\*(x<sup>(i)</sup>) + ξ<sup>(i)</sup> where the function h\*(·) is a quadratic function ξ<sup>(i)</sup> is the observation noise assumed to be generated from ~N(0,σ<sup>2</sup>).
- A test example (x,y) also has the same input-output relationship y = h\*(x) + ξ where ξ
   ~ N(0,σ<sup>2</sup>). It's impossible to predict the noise ξ, and therefore essentially our goal is to recover the function h\*(·)



The best fit linear model has large training and test errors





The issue cannot be mitigated with more training examples — even with a very large amount of, or even infinite training examples, the best fitted linear model is still inaccurate and fails to capture the structure of the data. Even if the noise is not present in the training data, the issue still occurs





a very (say, infinitely) large training dataset.

The linear model suffers from large bias, and underfits the data.

- The linear model family's inability to capture the structure in the data—linear models cannot represent the true quadratic function  $h^*$ , but not the lack of the data.
- Informally, we define the bias of a model to be the test error even if we were to fit it to



#### Next, we fit a 5th-degree polynomial to the data.



examples



#### The model learnt from the training set does not generalize well to other test



The bias of the 5-th degree polynomials is small—if we were to fit to an extremely large dataset, the resulting model would be close to a quadratic function and be accurate (because the family of 5-th degree polynomials) contains all the quadratic functions)

test error, called variance of a model fitting procedure



there is a large risk that we're fitting patterns in the data that happened to be present in our small, finite training set, but that do not reflect the wider pattern of the relationship between x and y. These "spurious" patterns in the training set are (mostly) due to the observation noise  $\xi(i)$ , and fitting these spurious patters results in a model with large test error.

#### The failure of fitting 5-th degree polynomials can be captured by another component of the

The **variance** can be intuitively characterized by the amount of variations across models learnt on multiple different training datasets (drawn from the same underlying distribution).

#### Bias and variance tradeoff

If our model is too "simple" (has very few parameters), then it may have large bias (but small variance), and it typically may suffer from underfitting.

large variance (but have smaller bias), and thus overfitting.

- If it is too "complex" and has very many parameters, then it may suffer from

# A mathematical decomposition of error (for regression)

$$MSE(x) = \sigma^{2} + \mathbb{E}[(h^{*}(x) - h)]$$
$$= \sigma^{2} + (h^{*}(x) - h)$$
$$= \sigma^{2} + (h^{*}(x))$$
unavoidable

- havg can be thought of as the best possible model learned even with infinite data
- The bias captures the part of the error that are introduced due to the lack of expressivity of the model. It is not due to the lack of data



- tradeoff
- can vary the size of the model (e.g., the number of features)
- model complexity and prevents overfitting

• Overftting is typically a result of using too complex models, and we need to choose a proper model complexity to achieve the optimal bias-variance

• When the model complexity is measured by the number of parameters, we

• However, the correct, informative complexity measure of the models can be a function of the parameters (e.g.,  $\ell_2$  norm of the parameters), which may not necessarily depend on the number of parameters. Then we use regularization

**Regularization** is an important technique in machine learning that controls the

denoted by  $R(\theta)$  here, to the training loss/cost function:

- small loss  $J(\theta)$  and have a small model complexity (a small  $R(\theta)$ ).
- bias.)

Regularization typically involves adding an additional term, called a regularizer and

 $J_{\lambda}(\theta) = J(\theta) + \lambda R(\theta)$ 

•  $J_{\lambda}$  is often called the regularized loss, and  $\lambda \ge 0$  is called the regularization parameter.

• The regularizer  $R(\theta)$  is typically chosen to be some measure of the complexity of the model  $\theta$ . Thus, when using the regularized loss, we aim to find a model that both fit the data (a

• When  $\lambda$  is a sufficiently small positive number, minimizing the regularized loss is effectively minimizing the original loss with the regularizer as the tie-breaker. When the regularizer is extremely large, then the original loss is not effective (and likely the model will have a large

- The most commonly used regularization is perhaps  $\ell_2$  regularization,  $R(\theta) = \frac{1}{2} \|\theta\|_2^2$ where
- standard gradient

$$\begin{aligned} \theta &\leftarrow \theta - \eta \nabla J_{\lambda}(\theta) = \theta - \eta \lambda \theta - \eta \nabla J(\theta) \\ &= \underbrace{(1 - \lambda \eta)\theta}_{\text{decominants}} - \eta \nabla J(\theta) \end{aligned}$$

decaying weights

• It encourages the optimizer to find a model with small  $\ell_2$  norm. In deep learning, it's oftentimes referred to as weight decay, because gradient descent with learning rate  $\eta$  on the regularized loss  $R_{\lambda}(\theta)$  is equivalent to shrinking/decaying  $\theta$  by a scalar factor of 1 –  $\eta\lambda$  and then applying the

## Regularization (sparsity)

- continuous surrogate the norm-1 (also called LASSO)

 $R(\theta)$ 

norm 2

 Besides encouraging simpler models, regularization can also impose inductive biases or structures on the model parameters. For example, suppose we had a prior belief that the number of non-zeros in the ground-truth model parameters is small, — which is oftentimes called sparsity of the model –, we can impose a regularization on the number of non-zeros in  $\theta$ 

• The sparsity of the parameters is not a continuous function of the parameters, and thus we cannot optimize it with (stochastic) gradient descent. A common relaxation is to use as a

$$= \|\theta\|_1$$

• Norms 1 and 2 are the most commonly used regularizers for linear models. In deep learning

## Riesgos sociales del AA (LLM)

#### GPTs are GPTs: An Early Look at the Labor Market Impact Potential of Large Language Models

Tyna Eloundou<sup>1</sup>, Sam Manning<sup>1,2</sup>, Pamela Mishkin\*<sup>1</sup>, and Daniel Rock<sup>3</sup>

<sup>1</sup>OpenAI <sup>2</sup>OpenResearch <sup>3</sup>University of Pennsylvania

March 20, 2023

- roles heavily reliant on science and critical thinking skills show a negative correlation with exposure, while programming and writing skills are positively associated with LLM exposure.
- information processing industries (4-digit NAICS) exhibit high exposure, while manufacturing, agriculture, and mining demonstrate lower exposure.

Job Zone	Preparation Required	Education Required	Example Occupations	Median Income	Tot Emp (000)	H a	N C
1	None or little (0-3 months)	High school diploma or GED (optional)	Food preparation workers, dishwashers, floor sanders	\$30,230	13,100	3.71	3.8
2	Some (3-12 months)	High school diploma	Orderlies, customer service representatives, tellers	\$38,215	73,962	7.03	11
3	Medium (1-2 years)	Vocational school, on-the-job training, or associate's degree	Electricians, barbers, medical assistants	54,815	37,881	11.28	13
4	Considerable (2-4 years)	Bachelor's degree	Database administrators, graphic designers, cost estimators	\$77,345	56,833	22.68	17
5	Extensive (4+ years)	Master's degree or higher	Pharmacists, lawyers, astronomers	\$81,980	21,221	22.81	13

Table 6: Exposure to GPTs by Job Zone



#### **Repercusiones IA**

- Disminuye la Intermediación
- los bancos, operadoras móviles, eléctricas, seguros ...
- Consumo energético

 Organizaciones (debido a la IA) transforman en plataformas gestionas por algoritmos que mueven una enorme cantidad de datos. Por ejemplo

## ¿Por dónde puede ir la solución?

- Humanidad: creatividad extra
- Laboral. Avanzar en
  - Intermediación 0
  - Automatización (plataformas) 0
  - **Optimización** de 0
    - procesos más sostenibles 0
    - la vida (cuidados) 0

## ¿Y la universidad?

- Cómo enseñamos?
- Para qué profesiones preparamos a nuestros estudiantes?
- Qué destrezas deberán tener?

## ¿Y ahora qué?

- Humildad género humano
- Humildad metodológica







When I think of existential risks to large parts of humanity: \* The next pandemic

\* Climate change  $\rightarrow$  massive depopulation

\* Another asteroid

thrive the next 1000 years, lets make AI go faster, not slower.

6:33 PM · May 30, 2023 · 1.1M Views

824 Retweets 175 Quotes 4,110 Likes 205 Bookmarks

#### ¿Y ahora qué?

...

- Al will be a key part of our solution. So if you want humanity to survive &





In large part because it doesn't exist yet.

level), discussing how to make it safe is premature.

#### ¿Y ahora qué?

- Super-human AI is nowhere near the top of the list of existential risks.
- Until we have a basic design for even dog-level AI (let alone human

...

