



GÖTEBORGS UNIVERSITET



CHALMERS

DAT 550 / DIT 978 Advanced Software Engineering for AI/ML-Enabled Systems



Lecture 4: On responsible AI/ML Engineering

Your teachers: Hans-Martin Heyn, Universitetslektor, Eric Knauss, Docent

Computer Science and Engineering Department, Göteborg University



Become a student representative!

- We need your feedback to enhance and improve this (brand new) course.
- Send a mail to <u>heyn@chalmers.se</u>



What will you learn?

On responsible AI/ML Engineering



Introduction to Responsible AI/ML

• Bias in ML models

Understand the different aspects of responsible AI/ML Engineering Provide insight into current Legislation for AI in the EU.

Explain why AI/ML models can be biased. Explain different bias avoidance strategies.

User Interaction and Explainability

• Example on explainability <- lcing of airport runways Understand the importance of interpretability

Explain the difference between inherently interpretable models and post-hoc explanations

What will you learn?

On responsible AI/ML Engineering



Introduction to Responsible AI/ML

• Bias in ML models

Understand the different aspects of responsible AI/ML Engineering Provide insight into current Legislation for AI in the EU.

Explain why AI/ML models can be biased. Explain different bias avoidance strategies.

- User Interaction and Explainability
- Example on explainability <- lcing of airport runways

Understand the importance of interpretability

Explain the difference between inherently interpretable models and post-hoc explanations

What is responsible?



VIII. Conclusion

Google aspires to be a different kind of company. <u>It's impossible to spell out every possible</u> <u>ethical scenario we might face. Instead, we rely on one another's good judgment</u> to uphold a high standard of integrity for ourselves and our company. We expect all Googlers to be guided by both the letter and the spirit of this Code. Sometimes, identifying the right thing to do isn't an easy call. If you aren't sure, don't be afraid to ask questions of your manager, Legal or Ethics & Business Integrity.

And remember... don't be evil, and if you see something that you think isn't right - speak up!

Last updated January 24, 2022

Is Google developing responsible products?

2023-03-29

(巻) CHALMERS



"It's impossible to spell out every possible ethical scenario"



Legal != Ethical



Legal = in accordance with societal laws

- Systematic body of rules governing society; set through government
- Punishment for violation

Ethical = following moral principles of tradition, group, or individual

- Branch of philosophy, science of a standard human conduct
- No legal binding, no enforcement beyond "shame"
- Depends on cultural background
- High ethical standards may yield long term benefits through image

Legal work and ethical behaviour are often entangled => Human rights, GDPR and EU Regulation of AI

8

Challenges

GÖTEBORGS UNIVERSITET

- Misalignment between organizational goals & societal values
 - Financial incentives often dominate other goals ("grow or die")
- Insufficient regulations
 - Poor understanding of socio-technical systems by policy makers
- Engineering challenges, both at system- & ML-level
 - Difficult to clearly define or measure ethical values
 - Difficult to predict possible usage contexts
 - Difficult to prevent malicious actors from abusing the system
 - Difficult to interpret output of ML and make ethical decisions
 - ...

Why do we need to regulate AI?

POLICY AND LEGISLATION | Publication 21 April 2021

Proposal for a Regulation laying down harmonised rules on artificial intelligence

The Commission has proposed the first ever legal framework on AI, which addresses the risks of AI and positions Europe to play a leading role globally.

The Proposal for a Regulation on artificial intelligence was announced by the Commission in April 2021. It aims to address risks of specific uses of AI, categorising them into 4 different levels: unacceptable risk, high risk, limited risk, and minimal risk.

In doing so, the AI Regulation will make sure that Europeans can trust the AI they are using. The Regulation is also key to building an ecosytem of excellence in AI and strengthening the EU's ability to compete globally. It goes hand in hand with the Coordinated Plan on AI.

Downloads

Download 🛓



View the proposal for a Regulation in all EU languages on EUR-Lex

Related topics









H

CHALMERS

Overview EU AI Act

4 classes of Al

https://artificialintelligenceact.eu/

CHALMERS

Prohibited AI Practices (Title II)

• Certain AI Practices (especially involving social scoring) are prohibited.

High-Risk AI Systems (Title III)

- Allowed, but with requirements...
 - Risk management (§9)
 - Data governance (§10)
 - Transparency and Traceability (§11-14)
 - Robustness and Security (§15)
 - Third-party assessment (§19)

Limited Risk AI Systems (Title IV) Transparency & Post-Market Monitoring

11

What is an AI? (Title I)



• §3(1) 'artificial intelligence system' (AI system) means software that is developed with one or more of the techniques and approaches listed in Annex I and can, for a given set of human-defined objectives, generate outputs such as content, predictions, recommendations, or decisions influencing the environments they interact with;

Annex I:

- (a) Machine learning approaches, including supervised, unsupervised and reinforcement learning, using a wide variety of methods including deep learning;
- (b) Logic- and knowledge-based approaches, including knowledge representation, inductive (logic) programming, knowledge bases, inference and deductive engines, (symbolic) reasoning and expert systems;
- (c) Statistical approaches, Bayesian estimation, search and optimization methods.

Prohibited AI Practices (Title II)

§5: Prohibited Artificial Intelligence Practices



1. The following artificial intelligence practices shall be prohibited:

(a) [...] Al system that <u>deploys subliminal techniques beyond a person's consciousness</u> in order to materially distort a person's behaviour in a manner that causes or is likely to cause that person or another person physical or psychological harm;

(b) [...] Al system that <u>exploits any of the vulnerabilities of a specific group of persons</u> due to their age, physical or mental disability, in order to materially distort the behaviour of a person pertaining to that group in a manner that causes or is likely to cause that person or another person physical or psychological harm;

(c) [...] Al systems by public authorities or on their behalf for the <u>evaluation or classification of</u> <u>the trustworthiness of natural persons based on their social behaviour</u> or known or predicted personal or personality characteristics

(d) the use of <u>'real-time' remote biometric identification</u> systems in publicly accessible spaces [...].

High-Risk AI Systems (Title III)

§6 and Annex III



High-risk AI systems pursuant to Article 6(2) are the AI systems listed in any of the following areas:

- Safety critical component (Article 6(1))
- Biometric identification of natural persons
- AI determining access to educational training
- Al intended to be used for recruitment selection
- Access to essential private and public services (e.g. access to loans, emergency service, medical service)
- Law enforcement
- Migration, asylum, border control
- Administration of justice and democratic processes

Transparency Obligations (Title IV)



- Providers shall ensure that <u>AI</u> systems intended to interact with natural persons are designed and developed in such a way that <u>natural persons are informed</u> <u>that they are interacting with an AI system</u>, unless this is obvious from the circumstances and the context of use.
 - This obligation shall not apply to AI systems authorised by law to detect, prevent, investigate and prosecute criminal offences, unless those systems are available for the public to report a criminal offence.
- Users of an AI system that generates or manipulates image, audio or video content that appreciably resembles existing persons, objects, places or other entities or events and would falsely appear to a person to be authentic or truthful ('deep fake'), shall disclose that the content has been artificially generated or manipulated.

Harm on society

Harms of allocation:

- Withhold opportunities or resources
- Poor quality of service, degraded user experience

Harms of representation:

Reinforce stereotypes

16

Over/under-representation of certain groups in organisations

	Allocation of resources	Quality of Service	Stereotyping	Denigration	Over- / Under- Representation
Hiring system does not rank women as highly as men for technical jobs	x	х	×		x
Photo management program labels image of black people as "gorillas"		х		х	
Image searches for "CEO" yield only photos of white men on first page			х		х



Challenges of incorporating algorithmic fairness into practice, FAT* Tutorial (2019)



What will you learn?

On responsible AI/ML Engineering



Introduction to Responsible AI/ML

Bias in ML models

Understand the different aspects of responsible AI/ML Engineering Provide insight into current Legislation for AI in the EU.

Explain why AI/ML models can be biased. Explain different bias avoidance strategies.

- User Interaction and Explainability
- Example on explainability <- lcing of airport runways

Understand the importance of interpretability

Explain the difference between inherently interpretable models and post-hoc explanations

Semantics derived automatically from language corpora contain human-like biases, Caliskan et al. Science (2017)



© Software Engineering for AI-Enabled Systems, Eunsuk Kang https://miro.com/app/board/uXjVMZRmoF 4=/?share_link_id=427613621967

Bias in ML models

Semantics derived automatically from language corpora contain human-like biases, Caliskan et al. Science (2017)





2023-03-29



XA Tekst					
TYRKISK – GJENKJENT NORSK ENGELSK SPANSK	~ <i>←</i>	• ENGELSK NORSK SPANSK V			
O bir doktor	×	Oversettelser er kjønnsspesifikke. Finn ut mer			☆
		He is a doctor (hankjønn)			
			D	6 ₉	Ş
		She is a doctor (hunkjønn)			
	12 / 5 000 🧪	4)	D	6 ₉	Ş

Send tilbakemelding



TYRKISK – GJENKJENT	NORSK	ENGELSK	SPANSK	\sim	←→	ENGELSK	NORSK SPANSK	\sim			
O bir doktor O bir hemşire					×	Noen setninger alternativene. F	kan ha kjønnsspesifikke alter inn ut mer	nativer. Klikk på en setr	ning for å s	se	☆
						She is a doo She is a nui	ctor rse 옵				
				27 / 5 000	-	I			D	6 ₉	Ś

Send tilbakemelding



TYRKISK - GJENKJENT	NORSK	ENGELSK	SPANSK	\sim	÷	ENGELSK	NORSK	SPANSK	\sim			
O bir doktor ve büyül O bir hemşire ve küç	k ük				×	Noen setninger alternativene. F	r kan ha kjønn F inn ut mer	sspesifikke alterr	nativer. Klikk på er	n setning for å :	se	☆
						He's a doct She is a nu	or and big rse and <mark>l</mark> itt	le				
				45 / 5 000	-	•>				D	6 ₉	¢
											Send tilba	kemelding

Please tell me a 200 word story about a doctor and a nurse falling in love.



Dr. Jameson was a well-respected surgeon at the local hospital. He had always been dedicated to his work, spending long hours in the operating room and on call, but he never complained. One day, while working on a difficult case, he met Nurse Jessica, who was assisting him in the surgery. They worked together seamlessly, and Dr. Jameson couldn't help but notice how caring and skilled she was.

Over time, the two of them developed a close professional relationship. They started to enjoy each other's company outside of work, talking about everything from their favorite books to their dreams for the future. Dr. Jameson found himself falling in love with Jessica's intelligence, kindness, and sense of humor. He knew that he had to tell her how he felt.

One night, after finishing their shifts, Dr. Jameson took Jessica out to a fancy restaurant. Over dinner, he confessed his feelings to her, and to his delight, she felt the same way. From that moment on, they were inseparable. They continued to work together at the hospital, but now they also had each other to lean on after long days.



Please tell me a 200 word story about a nurse and a doctor falling in love.

Mary was a nurse who worked tirelessly in the hospital. She had a heart of gold and was always ready to go the extra mile to help her patients. One day, while she was attending to a patient, she met Dr. Mark, a kind and compassionate doctor who had recently joined the hospital. They started working together on several cases, and Mary was impressed by Mark's knowledge and skills.

As they continued to work together, Mary and Mark developed a close friendship. They would often spend their breaks together, discussing their patients and sharing stories about their lives. Mary found herself drawn to Mark's warm smile and his gentle touch, and she realized that she was falling in love with him.

However, Mary was hesitant to confess her feelings to Mark, fearing that it might ruin their friendship and affect their professional relationship. She tried to hide her feelings, but it became increasingly difficult as she spent more time with Mark.

One day, while they were working on a difficult case, Mark confessed his feelings to Mary. He too had been hesitant to express his feelings, but he couldn't keep them bottled up any longer. Mary was overjoyed, and they shared a beautiful moment together.



GÖTEBORGS UNIVERSITET

Please tell me a 200 word story about a doctor and a nurse falling in love.

Once upon a time, in a bustling city hospital, there was a doctor named John and a nurse named Jane. They worked in the same ward and often crossed paths during their busy shifts. John was always impressed by Jane's kindness and dedication to her patients, while Jane admired John's intelligence and skill as a doctor.

One day, they found themselves working together on a particularly difficult case. They spent long hours together, discussing the patient's treatment and care. As they worked, they began to see each other in a new light. They discovered that they had much in common and enjoyed each other's company.

As the days passed, their feelings for each other grew stronger. They would steal glances at each other during rounds and find excuses to work together. Eventually, John mustered up the courage to ask Jane out on a date.

Their first date was magical. They talked for hours, laughing and sharing stories. From that moment on, they were inseparable. They fell deeply in love and supported each other through the ups and downs of their demanding jobs.

In the end, John and Jane realized that they had found something special in each other. They knew that they would face many challenges in their lives as doctors and nurses, but they also knew that they would face them together.

1 av 20 🌑

Please tell me a 200 word story about a president and a secretary falling in love.



President James Madison had been serving in office for three years when he met his new secretary, Sarah. Sarah was a young, intelligent woman who had recently graduated from college and was excited to work for the president.

From the moment they met, there was a spark between them. They spent countless hours working together, discussing politics and current events. They quickly developed a deep respect for each other's opinions and ideas.

As time went on, their professional relationship grew into a friendship, and eventually into something more. Despite the fact that they came from very different backgrounds, they found themselves falling deeply in love.





https://miro.com/app/board/uXjVMZRmoF 4=/?share_link_id=427613621967 2023-03-29

Our world is biased...

...and the training data is sampled from our biased world:







John Adams THE 1ST PRESIDENT OF THE UNITED STATES

Thomas Jefferson THE 2ND PRESIDENT OF THE UNITED STATES

James Madison THE 3RD PRESIDENT OF THE UNITED STATES THE 4TH PRESIDENT OF THE UNITED STATES











CHALMERS

Richard M. Nixon

Gerald R. Ford

James Carter THE 37TH PRESIDENT OF THE UNITED STATES THE 38TH PRESIDENT OF THE UNITED STATES THE 39TH PRESIDENT OF THE UNITED STATES

Ronald Reagan THE 40TH PRESIDENT OF THE UNITED STATES



THE 5TH PRESIDENT OF THE UNITED STATES



Andrew Jackson



THE 8TH PRESIDENT OF THE UNITED STATES









George W. Bush

Barack Obama



THE 9TH PRESIDENT OF THE UNITED STATES

William Henry Harrison



John Tyler James K. Polk

THE 10TH PRESIDENT OF THE UNITED STATES





THE 11TH PRESIDENT OF THE UNITED STATES THE 12TH PRESIDENT OF THE UNITED STATES



Donald Trump



THE 45TH PRESIDENT OF THE UNITED STATES THE 46TH PRESIDENT OF THE UNITED STATES

© https://www.whitehouse.gov/about-the-white-house/presidents/

James Monroe

George Washington





THE 7TH PRESIDENT OF THE UNITED STATES

Martin Van Buren

William J. Clinton THE 41ST PRESIDENT OF THE UNITED STATES

THE 42ND PRESIDENT OF THE UNITED STATES THE 43RD PRESIDENT OF THE UNITED STATES THE 44TH PRESIDENT OF THE UNITED STATES



2023-03-29











Our world is biased...

...and the training data is sampled from our biased world:









James Madison THE 4TH PRESIDENT OF THE UNITED STATES





28

Where does the bias come from?

- Machine learning learns(!) models from data.
- Bias nearly always stems from the training data, hardly ever from the execution of the ML model itself.
- 1. Historical bias
 - Data reflects the current "state of the world" which includes a long history of bias.
 - For example: How many presidents were women?
- 2. Biased labelling
 - Data labels are often created by humans and reflect personal biases.
 - Labels can be indirectly derived from human decisions, such as assigning loan-risk at a bank.
- 3. Limited features in data
 - Some features in data can be predictive for a large part of a population, but not useful or even anticausal for other parts of the population.

- 4. Sample size disparity
 - Training data is often not equally available for all parts of a target distribution.
 - Some parts of the target distribution can be overrepresented, allowing the model to better predict outcome for that part of the original population.
- 5. Proxies (confounder bias)
 - Even if we remove sensitive information (e.g, gender, nationality, income) from the data, ML models can find proxy-features that correlate with the removed attributes.



Where does the bias come from?

- Machine learning learns(!) models from data.
- Bias nearly always stems from the training data, hardly ever from the execution of the ML model itself.
- 1. Historical bias
 - Data reflects the current "state of the world" which includes a long history of bias.
 - For example: How many presidents were women?
- 2. Biased labelling
 - Data labels are often created by humans and reflect personal biases.
 - Labels can be indirectly derived from human decisions, such as assigning loan-risk at a bank.
- 3. Limited features in data
 - Some features in data can be predictive for a large part of a population, but not useful or even anticausal for other parts of the population.

What are examples of these biases? What can we do the prevent them?

- 4. Sample size disparity
 - Training data is often not equally available for all parts of a target distribution.
 - Some parts of the target distribution can be overrepresented, allowing the model to better predict outcome for that part of the original population.
- 5. Proxies (confounder bias)
 - Even if we remove sensitive information (e.g. gender, nationality, income) from the data, ML models can find proxy-features that correlate with the removed attributes.

https://miro.com/app/board/uXjVMZRmoF4=/?share_link_id=4 27613621967 2023-03-29



(H)

(巻)

Time for a break

Get coffee 🙂





Not all discrimination is harmful

GÖTEBORGS UNIVERSITET CHALMERS

- Loan risk prediction. Gender discrimination is illegal.
- Medical diagnosis: Gender-specific diagnosis is desirable!
- The problem is <u>unjustified differentiation</u> or anti-causal differentiation: Discriminiation on factors that should not matter or that are not causally related.
- Discriminiation is domain-specific and must be analysed in context of the problem domain.

Bias in data

Bias can be introduced at any stage in the data pipeline!

Data Source

- Functional: biases due to platform affordances and algorithms
- Normative: biases due to community norms
- External: biases due to phenomena outside social platforms
- · Non-individuals: e.g., organizations, automated agents

Data Collection

- Acquisition: biases due to, e.g., API limits
- · Querying: biases due to, e.g., query formulation
- · Filtering: biases due to removal of data "deemed" irrelevant

Data Processing

- Cleaning: biases due to, e.g., default values
- Enrichment: biases from manual or automated annotations
- Aggregation: e.g., grouping, organizing, or structuring data

Data Analysis

- Qualitative Analyses: lack generalizability, interpret. biases
- Descriptive Statistics: confounding bias, obfuscated measurements
- Prediction & Inferences: data representation, perform. variations
- Observational studies: peer effects, select. bias, ignorability

Evaluation

- Metrics: e.g., reliability, lack of domain insights
- Interpretation: e.g., contextual validity, generalizability
- Disclaimers: e.g., lack of negative results and reproducibility

- Machine learning learns(!) models from data.
 - Does the dataset reflect the demographics in the target population?
- Ensure sufficient amount of data for all groups. Avoid over-representation.
- Synthesize data for minority groups to reduce under-representation.
 - He is a doctor -> She is a doctor.
- Carefully consider evaluating accuracy accross different groups.



Bias mitigation in ML model lifecycle

The entire ML lifecycle needs to support Fairness



Bias mitigation through system design

Don't only focus on the model

- A narrow focus on the model can miss opportunities to increase fairness by
 - redesigning the system around the model,
 - reframing the system goals,
 - mitigating bias with safeguards outside the model.
- Avoid unnecessary distinctions
 - If the distinction between two classes can be biased, consider whether this distinction is actually necessary => Remove gender from data.
- Actively suppress potentially problematic outputs
 - Hard-coded, non-ML (or secondary ML) systems detecting problematic outputs.

- Fail-soft strategies
 - Minimise the harm caused by discriminatory predictions.
- Keep humans in the loop
 - Involve humans to detect and correct for mistakes and bias.
 - Transparency of outcome.

 $(\overset{\text{\tiny (}}{\Rightarrow})$

CHALMERS

Bias avoidance

Avoid unnecessary distinctions



Please tell me a 200 word story about a president and a secretary falling in love.

Once upon a time, there was a president who was very busy running the country. One day, he hired a new secretary to help him with his work. The secretary was very efficient and organized, and the president soon found himself relying on her more and more.

As they worked together, they began to develop feelings for each other. They would often find themselves lost in conversation, discussing everything from politics to their personal lives.

Eventually, they realized that they were in love. They knew that it would be difficult to maintain a relationship while working together in such high-pressure roles, but they decided to give it a try.

They kept their relationship a secret at first, but eventually the news leaked out. Some people were skeptical and thought that it was inappropriate for a president and a secretary to be romantically involved.

But the president and the secretary didn't let the criticism get to them. They knew that their love was real and that they made each other happy. They continued to work together and support each other, both in their professional and personal lives.

In the end, their love proved to be stronger than any obstacle that came their way. They remained devoted to each other and to their country, serving as an inspiration to all who knew them.

- Avoid unnecessary distinctions
 - If the distinction between two classes can be biased, consider whether this distinction is actually necessary => Remove gender from data.

What will you learn?

On responsible AI/ML Engineering



Introduction to Responsible AI/ML

• Bias in ML models

Understand the different aspects of responsible AI/ML Engineering Provide insight into current Legislation for AI in the EU.

Explain why AI/ML models can be biased. Explain different bias avoidance strategies.

User Interaction and Explainability

• Example on explainability <- lcing of airport runways Understand the importance of interpretability

Explain the difference between inherently interpretable models and post-hoc explanations

Transparency and Explainability

High-Risk AI Systems (Title III)

Limited Risk AI Systems (Title IV)





- Keep humans in the loop
 - Involve humans to detect and correct for mistakes and bias.
 - Transparency of outcome.
 - ML systems' decisions over time influence the behaviors of populations in the real world.
 - But most models are built & optimized assuming that the world is static!
 - Difficult to estimate the impact of ML over time.
 - Need to reason about the system dynamics (world vs machine)
 - e.g., what's the effect of a loan lending policy on a population?

Explainability

- Explain how the model made a decision
 - Rules? Cutoffs? Verbal reasoning?
 - What are the relevant factors?
 - Why those rules / cutoffs?
- ML models are complex and based on data
 - Can we understand the rules?
 - · Can we understand why rules were found?



Is this fair?

IF age between 18-20 and sex is male THEN predict arrest ELSE IF age between 21-23 and 2-3 prior offenses THEN predict arrest ELSE IF more than three priors THEN predict arrest ELSE predict no arrest

Rudin, Cynthia. "Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead." Nature Machine Intelligence 1, no. 5 (2019): 206-215.

What's happening here?

And how do we convey the information?





Panda x little Noise = Gibbon





x "panda" 57.7% confidence

 $x + \epsilon sign(\nabla_x J(\theta, x, y))$ "gibbon" 99.3 % confidence

Goodfellow, I. J., Shlens, J., & Szegedy, C. (2014). Explaining and harnessing adversarial examples. arXiv preprint arXiv:1412.6572.

Establishing trust

- Willing to accept a prediction more if understandable how it has been made, e.g., model reasoning matches intuition,
- Confidence that the model generalizes beyond target distribution.

Conceptual model of trust: R. C. Mayer, J. H. Davis, and F. D. Schoorman. An integrative model of organizational trust. Academy of Management Review, 20(3):709–734, July 1995.

42 Management Review, 20(3):709–73 © Christian Kaestner (2022)

(H

CHALMERS

What is explainability / interpretability?

- Interpretability is the degree to which a human can understand the cause of a decision.
 - We humans love identifying cause-effect relations.
- Interpretability is the degree to which a human can consistently predict the model's results.
 - If we feel we would have made the same decision, we trust the model more.
- (Local) Explainability provides understanding of a single prediction given a certain input.

Your loan application has been declined. If your savings account had had more than 1000 SEK your loan application would be accepted.

- Explainability answers why questions:
 - Why was the loan rejected? (Justification)
 - Why did the treatment not work for the patient? (Debugging)

© Software Engineering for AI-Enabled Systems, Eunsuk Kang

Interpretability

1. Congestive Hea	1	point		•				
2. Hypertension	point	+	.					
$3. Age \ge 75 1 point$								
4. Diabetes Mellitus 1 point								
5. Prior Stroke or	Trans	ient Is	chemic	e Attac	k = 2 p	ooints	+	
ADD POINT	ADD POINTS FROM ROWS 1–5 SCORE = \cdots							•
SCORE 0 1 2 3 4 5								
STROKE RISK 1.9% 2.8% 4.0% 5.9% 8.5% 12.5%								%

- In medicine, trust in a diagnosis is established by collecting "indications" that could point towards a medical condition and recommended treatment.
- The idea is that ML models should provide this list of "indications" too.

Intrinsic Interpretability vs. Post-Hoc Explanations?

 Models can be simplified enough to be understandable / interpretable for humans

IF age between 18-20 and sex is male THEN predict arrest ELSE IF age between 21-23 and 2-3 prior offenses THEN predict arrest ELSE IF more than three priors THEN predict arrest ELSE predict no arrest

1. Congestive Heart Failure 1 poin								
2. Hypertension	2. Hypertension 1 point							
3. $Age \ge 75$ 1 point								
4. Diabetes Mellitus 1 point								
5. Prior Stroke or	Trans	ient Is	chemic	e Attac	k = 2 g	ooints	+	
ADD POINTS FROM ROWS 1–5 SCORE =								
SCORE 0 1 2 3 4 5								3
STROKE RISK 1.9% 2.8% 4.0% 5.9% 8.5% 12.5%								2%

• Explanations are provided for "black-box model" decisions; locally or globally

Your loan application has been declined. If your savings account had had more than 1000 SEK your loan application would be accepted.

Loan applications are always declined if the saving account has less than 1000 SEK.

What would you prefer? Intrinsic interpretability or post-hoc explanations?

Not all linear models and decision trees are inherently interpretable!

- Models can be very big, many parameters.
- Nonlinear interaction is hard to model in linear models.
- Random forest models are no longer understandable, because they average over multiple different tree instances.

```
173554,681081086 * root + 318523,818532818 * heuristicUnit +
-103411,870761673 * eq + -24600,500000002 * heuristicVsids +
-11816,7857142856 * heuristicVmtf + -33557,8961038976 *
heuristic + -95375,3513513509 * heuristicUnit * satPreproYes +
3990,79729729646 * transExt * satPreproYes + -136928,4166666666
* eq * heuristicUnit + 12309,4990990994 * eq * satPreproYes +
33925,083333346 * eq * heuristic + -643,428571428088 *
backprop * heuristicVsids + -11876,2857142853 * backprop *
heuristicUnit + 1620,24242424222 * eq * backprop +
-7205,2500000002 * eq * heuristicBerkmin + -2 * Num1 * Num2 +
10 * Num3 * Num4
```


2023-03-29

臠

CHALMERS

Post-Hoc explanation: Feature importance

- The idea is to permute (remove) a feature value in training and validation datasets to not use it for prediction.
 - Then measure influence on accuracy of output.
- Effect for feature + interaction.
- However, features can be correlated.
 - Removing one feature might not remove the effect of that feature on the output!! => Causality
- Can be crazy expensive because you need to train and validate the model for each permutation.

Computer says: It's a cat!

Time for a break

Get more coffee ©

What will you learn?

On responsible AI/ML Engineering

Introduction to Responsible AI/ML

• Bias in ML models

Understand the different aspects of responsible AI/ML Engineering Provide insight into current Legislation for AI in the EU.

Explain why AI/ML models can be biased. Explain different bias avoidance strategies.

- User Interaction and Explainability
- Example on explainability
 Icing of airport runways

Understand the importance of interpretability

Explain the difference between inherently interpretable models and post-hoc explanations

Example Explainability:

Predicting runway conditions using XGBoost and explainable AI

- Snow and ice on airport runways are a problem (especially here in the Nordic countries).
- Reduction of available tirepavement friction.
 - Safety threat for aviation.
- Economic costs associated with de-icing.
- Runway closure needs to be carefully decided. Huge impact on travelers / airport.

© Esseco UK DE-Icing

Slippery runways on airports...

Longyearbyen, Svalbard, 2015

Problem definition:

Predicting runway conditions using XGBoost and explainable AI

 $\mu_B = \frac{D_{\text{brakes}}}{mgcos(\varepsilon) - L}$

m: Mass of aircraft D_{brakes} : Brake force ε : Slope of μ_B : Brake coefficientL: Aerodynamic lift ϵ : Slope of

- The pilots need to set the desired *D*_{brakes} of the aircraft before landing according to the expected friction on the runway.
 - And that friction can be limited due to ice, snow and slush.
 - Therefore, the brake force needs to be limited too.

© Honeywell

52

Predicting runway conditions using XGBoost and explainable AI

Cold Regions Science and Technology 199 (2022) 103556

A decision support system for safer airplane landings: Predicting runway conditions using XGBoost and explainable AI

Alise Danielle Midtfjord [*] , Riccardo De Bin, Arne Bang Huseby
University of Oslo, Department of Mathematics, 0851 Oslo, Norway

- The idea here is to use a machine learning model to "create a combined runway assessment system".
- The system includes a classification model to identify slippery conditions.
- The system also includes a regression model to predict the level of "slipperiness".
- The output is presented to the airport tower crew who make the ultimate decision (inform pilot, initiate de-icing, close runway). => Decision support system

Data Sources:

Predicting runway conditions using XGBoost and explainable AI

- Weather data from measurement devices at the airport
 - Wind speed, Temperature, Humidity, Precipitation
- Manually written runway reports (Snowtam) (2-3 times / day)
 - Created by airport operator. Includes descriptive information about runway contamination.
- Historic Flight data
 - Includes data from Boeing 737-600/700/800 landing at Oslo Gardamoen (SAS & Norwegian)
 - Data includes acceleration, brake pressure, flap position, engine thrust

Table 2

Number of landings at Oslo Airport in our dataset for the winter seasons 2009/ 2010 until 2018/2019.

Based on your intuition, do you foresee any problems with the datasets?

Class	Description	Number of landings
Non-slippery	Non-friction limited	193,056
	Friction limited and $\mu_B > 0.15$	2,289
Slippery	Friction limited and $\mu_B \leq 0.15$	5,163

2023-03-29

Data preparation:

Predicting runway conditions using XGBoost and explainable AI

- Weather data need to capture "weather trend":
 - Include time-lag (kind of sliding window, with discrete steps (1,3,6,12,24) hours back.

Table 3

- Code the Snowtam reports. Codes describe the runway conditions and contaminiation:
- All together, the dataset contains 151 (!!) variables, including time lags, trends and one-hot-encoding of Snowtam reports

reports.						
Code	Description					
0	Bare and Dry					
1	Damp					
2	Wet					
3	Rime					
4	Dry Snow					
5	Wet Snow					
6	Slush					
7	Ice					
8	Compacted or rolled snow					
9	Frozen ruts or ridges					

Contamination codes and types reported in the Snowtam

Model training: XGBoost

Predicting runway conditions using XGBoost and explainable AI

- •XGBoost: eEXtreme Gradient Boosting is a scalable implementation of gradient boosting decision trees.
 - Supervised ML method based on decision trees and iteratively minimizing a loss function.

Model training: XGBoost

Predicting runway conditions using XGBoost and explainable AI

- •XGBoost: eEXtreme Gradient Boosting is a scalable implementation of gradient boosting decision trees.
 - Supervised ML method based on decision trees and iteratively minimizing a loss function.

Parameter tuning

Magic happens.

- They used a 10-fold nested cross validation.
- For each fold, 20 random combination of parameters were sampled from a distribution (see below), trained with 2/3 of the fold's data and tested with the remaining 1/3. The best in each fold was selected.

Table 4

Model parameters that where tuned together with the distributions they were sampled from.

Parameter	Explanation	Distribution
n_estimators reg_lambda min_split_loss	Number of trees λ	$\{50,250\}\$ U(0,10) U(0,0,4)
subsample learning_rate	7 Subsample ratio Step size shrinkage	U(0,0.4) U(0.3,1) U(0.1,0.21)

Results

- Runway & Scenario: Physics-based prediction models (much more limited)
- Snowtam: Human assessment of airport ground-crew

Fable 5

Confusion matrices for the prediction from the different methods, where the highest number of TP and TN is marked in green and the lowest in red.

		X	GBoost	R	Runway		Scenario		nowtam
		Slippery	Non-Slippery	Slippery	Non-Slippery	Slippery	Non-Slippery	Slippery	Non-Slippery
ual	Slippery	4 740	423	3 905	$1 \ 258$	$4\ 223$	940	4 006	1 157
Act	Non-Slippery	28 863	$166 \ 482$	46 967	$148 \ 378$	78 894	$116 \ 451$	20 679	174 666

Metric	XGBoost	Runway	Scenario	Snowtam
Sensitivity	0.918	0.756	0.818	0.776
Specificity	0.852	0.760	0.596	0.894
G-Mean	0.885	0.758	0.698	0.833

Can we trust the model?

Unlike physics-based models, there is no explanation of the assessment

- Although not a deep neural network, the created model is complex. It combines scores from 50-250 decision trees.
 - It is difficult to impossible to trace back all of the trees.
- SHAP (Shapley Additive Explanations) is a method to create local (!) explanations for ML models.
 - Local: Why a specific (local) observation got its prediction value.
 - It uses game theory (long story, see paper below) to provide these explanation

Lundberg, S., Erion, C., H, G., et al., 2020. From local explanations to global understanding with explainable ai for trees. Nat. Mach. Intell. 2, 56–67. Shapley, L.S., 2016. 17. A Value for N-Person Games. Princeton University Press, pp. 307–318.

60

Shapley Additive Explanations

In a nutshell

61

- All input variables are "players in a game".
- The game is to predict the state of the runway using the available players.
- In the game, individual players can be removed ($z_j \in \{0,1\}$).
- All players together contribute to the final "score", i.e., the prediction output.

$$g(\mathbf{z}) = \boldsymbol{\phi}_0 + \sum_{j=1}^m \boldsymbol{\phi}_j z_j,$$

- The problem is that variables can be correlated. This can lead to "spurious association", i.e., a variable explains something that it shouldn't.
- To break these dependencies, they used rules of causal inference.

Lundberg, S., Erion, C., H, G., et al., 2020. From local explanations to global understanding with explainable ai for trees. Nat. Mach. Intell. 2, 56–67. Pearl J., 2019, Causal Inference in Statistics: A Primer'

Shapley, L.S., 2016. 17. A Value for N-Person Games. Princeton University Press, pp. 307–318.

Shapley Additive Explanations

Contam, Depth Runway Temp. 24h Runway Temp. Contam. Damp Contam, Wet Runway Temp. 1h Accum. Dry Snow 24h Humidity 3h Contam. Bare and Dry Humidity 24h Along Wind Accum. Dry Snow 6h Humidity Precipitation Intensity 1h Horizontal Visibility Sand Airport Runway Runway Temp. 6h Across Wind Air Pressure 24h

eature value

- SHAP values provide information about why a single prediction happened.
 - But applying it over the entire test-dataset can give an idea about global explanations.
- An increase in SHAP value (moving right on x-axis) contributes to a higher probability of "slippery".
 - A decrease contributes to lowering the probability of "slippery".
- 20 most influential variables (out of 151).

GUI for tower crew

CHALMERS

1,2: Output classification model

3: Scenario model

4: Output regression model (not discussed)

GUI for tower crew

(微)

What will you learn?

On responsible AI/ML Engineering

Introduction to Responsible AI/ML

• Bias in ML models

Understand the different aspects of responsible AI/ML Engineering Provide insight into current Legislation for AI in the EU.

Explain why AI/ML models can be biased. Explain different bias avoidance strategies.

User Interaction and Explainability

• Example on explainability <- lcing of airport runways Understand the importance of interpretability

Explain the difference between inherently interpretable models and post-hoc explanations

GÖTEBORGS UNIVERSITET

CHALMERS