

A Five-Phase Framework for Fair Insurance: Reviewing Strategies for Digital Price Differentiation

RIJK MERCUUR, HU University of Applied Sciences Utrecht, Netherlands

SIEUWERT VAN OTTERLOO, HU University of Applied Sciences Utrecht, Netherlands

HUIB ALDEWERELD, HU University of Applied Sciences Utrecht, Netherlands

Insurers increasingly use machine learning to assess financial risk and determine personalized premiums with the aim of ensuring income and financial stability. Despite the recent focus on fairness, insurance companies and software developers struggle to bridge the gap between fairness principles and practical implementation. Justifying digital price differentiation in terms of both fairness and profit is a socio-technical problem: it requires an integration of organizational processes, ethical-legal considerations on indirect discrimination, and technical fairness metrics and mitigation techniques. The paper proposes a structured list of 33 strategies designed to help organizations navigate these challenges, derived from a survey of Dutch insurance professionals, a systematic review of academic literature, and expert evaluations. The strategies are organized into five phases: Understand, Determine, Adjust, Evaluate and Communicate, with a particular emphasis on aligning fairness principles with actuarial accuracy and compliance with legal standards. This work contributes to the literature by offering an overview of actionable strategies that go beyond fairness metrics, addressing both technical and social aspects of digital price differentiation. Practically, the strategy list supports insurance professionals — including data scientists, actuaries, auditors, compliance officers, and communication staff — by (1) providing a comprehensive overview of strategies to balance fairness and profitability in digital price differentiation, and (2) offering a framework to structure organizational processes and internal communication around this balance.

Keywords: strategies, insurance, pricing, price differentiation, price discrimination, fairness metrics, algorithmic bias, solidarity, discrimination, literature review, AI, machine learning

Reference Format:

Rijk Mercuur, Sieuwert van Otterloo, and Huib Aldewereld. 2025. A Five-Phase Framework for Fair Insurance: Reviewing Strategies for Digital Price Differentiation. In *Proceedings of Fourth European Workshop on Algorithmic Fairness (EWAF'25)*. Proceedings of Machine Learning Research, 21 pages.

1 Introduction

Financial institutions increasingly implement artificial intelligence (AI) to improve financial performance, yet this raises concerns about the fair treatment of clients [53]. As of 2024, 37% of Dutch financial institutions have adopted AI [9]. In particular, insurance companies and software providers in the insurance sector use AI to calculate individualized premiums based on the specific circumstances of the client, a practice known as price

Authors' Contact Information: Rijk Mercuur, HU University of Applied Sciences Utrecht, Utrecht, Netherlands, rijk.mercuur@hu.nl; Sieuwert van Otterloo, HU University of Applied Sciences Utrecht, Utrecht, Netherlands; Huib Aldewereld, HU University of Applied Sciences Utrecht, Utrecht, Netherlands, huib.aldewereld@hu.nl.

This paper is published under the Creative Commons Attribution-NonCommercial-NoDerivs 4.0 International (CC-BY-NC-ND 4.0) license. Authors reserve their rights to disseminate the work on their personal and corporate Web sites with the appropriate attribution.

EWAF'25, June 30–July 02, 2025, Eindhoven, NL

© 2025 Copyright held by the owner/author(s).

differentiation. Following Zuiderveen Borges [52], we refer to the use of machine learning to determine these individual premiums as *digital price differentiation*. Insurance is an essential service, and reasonable pricing that does not make insurance unaffordable for specific groups is therefore important. Consequently, insurance companies, regulators, legislators, and other stakeholders seek to ensure that any price differentiation does not lead to discrimination or unfairness [19, 28, 30, 52].

Balancing the need for accurately pricing risk while avoiding unfair bias and discrimination is challenging for insurance companies and their software developers. This challenge stems from a broader socio-technical dilemma in which technical solutions alone cannot resolve the ethical and strategic complexities of fairness in digital price differentiation. Four key factors contribute to this dilemma:

- (1) The significant risks to insurance companies if their pricing models are not accurate due to adverse selection [10]. If insurers' prices are too high, they will lose customers to the competition since many consumers use price comparison websites and are price-sensitive [10]. If prices are too low, insurers will attract more customers but will have to pay more out in expected damages than they will earn in premiums for each customer, effectively losing money.
- (2) A fundamental tension between fairness and maximizing financial returns. Insurers aim to refine digital price differentiation to reflect individual risks as precisely as possible. However, fairness considerations often require treating similar clients equitably, even when granular risk assessments suggest otherwise. For example, insurers may use ZIP codes to differentiate premiums based on neighborhood-based risk variations, increasing returns while potentially exacerbating socio-economic disparities. This inherent conflict requires organizations to navigate trade-offs between fairness principles and actuarial accuracy.
- (3) Ambiguity remains regarding how to mitigate indirect discrimination and translate fairness principles into organizational strategies. While direct discrimination—explicit differentiation based on protected characteristics such as gender or ethnicity—is relatively straightforward to identify and eliminate [32, 49], indirect discrimination is more complex. Seemingly neutral variables, such as ZIP codes or vehicle types, may serve as proxies for protected attributes, leading to disparate impacts [20, 25, 32, 35, 51]. The legal and ethical boundaries of such practices remain uncertain, making it difficult for insurers to establish clear and defensible fairness policies [14, 52].
- (4) Integrating scientific and technical methods for improving and measuring fairness into social processes presents significant challenges. The literature provides various fairness metrics and bias mitigation techniques (e.g., synthetic data augmentation) [36]. However, selecting appropriate metrics, defining target fairness thresholds, and implementing mitigation strategies require integration of ethical, legal, technical, and organisational aspects. Effective integration demands not only technical expertise but also organisational alignment on understanding, mitigating, evaluating and communicating fairness. Especially, since the insurance industry is regulated and insurance pricing experts need to present their pricing models and risk calculations to internal and external supervisors. The requirements and expectations of stakeholders put additional social and organisational constraints of what models or inputs can be used.

In short, addressing fairness in digital price differentiation requires more than just technical tools or laws — it demands a structured socio-technical approach that also involves legal checks and communication. To support

professionals in this balancing act, a concrete and ordered list of strategies is needed. The list presented below offers clear guidance on how insurers and software providers can evaluate and justify digital price differentiation.

This paper addresses the question:

"What can insurance companies and software providers for the insurance sector do to justify the use of digital price differentiation from the perspectives of both fairness and financial returns?"

Unlike optimization-focused approaches, which aim at maximizing outcomes, our focus is on justification—ensuring that digital price differentiation meets ethical, legal, and financial standards. We define a strategy as a series of actions aimed at upholding or achieving (organizational) values [42]. To answer our research question, we create a structured list of 33 strategies to help insurance companies and software providers justify digital price differentiation while balancing fairness and financial return (see Table 3-8).

The primary contribution of this paper is this structured strategy list, which was created through a survey of Dutch insurance practitioners, a systematic review of academic and industry papers, and evaluated with professionals and experts. Scientifically, this list is unique for its comprehensive scope, addressing both technical and social aspects, and for its focus on digital price differentiation, which is a regression problem, rather than a classification one. Our work goes beyond fairness metrics to offer practical strategies that fit into organizational processes. Practically, the strategy list supports insurance professionals — including data scientists, actuaries, auditors, compliance officers, policymakers, and public relations teams — by (1) providing a comprehensive overview of strategies to balance fairness and profitability in digital price differentiation, and (2) offering a framework to structure organizational processes and internal communication around this balance.

The remainder of the paper is structured as follows: Section 2 reviews related work on legislation, protocols, and technical metrics; Section 3 details the scoping of the literature review, the systematic selection of literature from academia and practice, and the construction of the strategy list; and Section 4 presents the resulting strategy list, followed by a conclusion and suggestions for future research. As described in Appendix A: a full interactive version of the strategy list including explanations per strategy can be found on <https://explainable-ai.nl/4d/>. Appendix B shows the methodology and result of the survey used for scoping the literature. Appendix C contains the full strategy list (i.e., Table 3-8).

2 Background & Related Work

This section reviews current work aimed at optimizing digital price differentiation for insurance companies and software developers, focusing on Dutch and European regulations. While regulations across European countries vary, they are largely similar due to overarching EU laws, such as the AI Act, GDPR, and the European Convention on Human Rights.

2.1 Dutch and European Union Legislation

The European Convention on Human Rights (Article 14) prohibits discrimination based on attributes such as sex, race, religion, and other status. This primarily addresses direct discrimination. The Convention does not specify whether indirect discrimination is forbidden, leaving room for interpretation by individual countries.

Dutch law, particularly Article 1 of the Dutch Constitution and the General Equal Treatment Act [35], forbids direct discrimination based on categories race, religion, gender, and sexual orientation (Article 1.1.1c) but not age.

The law also prohibits indirect discrimination, defined as practices that disproportionately affect protected groups (Article 1.1.1c), unless such practices have a legitimate aim, are necessary and proportionate (Article 2.2.1). This leaves insurers and regulators responsible for ensuring that their practices are justifiable under these terms.

At the EU level, laws such as the GDPR and the AI Act emphasize non-discrimination and fairness but leave specific interpretations of indirect discrimination to existing legal frameworks and public interpretations [15, 53]. The GDPR stresses fairness (Article 5) but lacks detailed criteria for measuring fairness or assessing proportionality and necessity. The AI Act, intended to address AI-related risks, attempts to protect human rights by requiring AI providers to examine data for possible biases that could lead to discrimination. However, the Act does not provide clear guidelines on what constitutes an adequate examination of such biases, thresholds, or how they should be measured. As noted by Deck et al. [15], "the AI Act lacks clear substantive standards for determining when unequal treatment is inadmissible".

Although the AI Act is particularly aimed at high-risk AI, such as people's health or work performance are assessed, similar principles apply to limited- and low-risk AI applications, property insurance. The European Insurance and Occupational Pensions Authority (EIOPA) highlights that fairness considerations must be incorporated into risk assessments, regardless of whether the AI is classified as high-risk. A 2025 consultation paper from EIOPA reinforces this, emphasizing the need for fairness governance and risk assessments for all AI applications in insurance [19].

2.2 Protocols & Technical Metrics

Current law does not specify how to measure fairness and how to determine what level of fairness is sufficient. Practitioners aim to bridge this gap by using comparisons with previously approved models applying guidelines or using checklists. One widely cited fairness guideline is the four-fifths rule [6], which states that the hiring rate of any protected group must be no less than 80% of the rate for non-protected groups. This rule, developed in the US in the 1970s, provides a clear and easy-to-implement benchmark but can be arbitrary when applied to multiple protected groups or complex scenarios. As Deck et al. [15] observes, "technical fairness metrics such as statistical parity or equalized odds offer an actionable approach to measure and mitigate 'bias'. However, it remains unanswered what kind of evidence would signal sufficient efforts of bias detection and correction" [15].

There are multiple examples of checklists designed to help practitioners with fairness decisions. Directly related to our aim is the socio-technical fairness checklist by Madaio et al. [34]. This checklist, co-designed to tackle organizational challenges related to fairness, focuses on aligning fairness assessments with the needs of practitioners. As Madaio et al. [34] states, other fairness checklists [1, 8, 12, 17, 23, 24, 28] are too broad, overly specific, or narrow in scope, often focusing solely on aspects like data collection. Madaio et al.'s approach is distinct because it emphasizes co-designing tools that match practitioners' needs, rather than providing rigid checklists. We extend their work by applying a similar framework to the specific context of insurance and digital price differentiation, focusing on actionable strategies and addressing the balance between fairness and financial returns, particularly through communication and justification.

Dutch ethical protocols, such as DEDA [46] and IAMA [45], emphasize the need to evaluate fairness based on human rights principles. While these frameworks offer useful structured approaches, they remain general and not sector-specific, lacking specific guidance on operationalizing fairness in practice. On an international scale, the

ISO 42001 standard for AI management identifies fairness as an essential component of AI risk assessment but does not provide concrete instructions for implementing fairness within organizational contexts.

In terms of technical fairness metrics, there are several tools available, such as statistical parity and equalized odds, which are designed to assess and mitigate bias in AI systems [6]. However, there is no consensus on what constitutes an adequate level of fairness or what evidence is needed to demonstrate sufficient efforts at bias mitigation. Therefore, technical metrics need to be part of a broader, structured process that helps practitioners justify digital price differentiation in a transparent and fair manner.

3 Literature Selection

3.1 Type of Literature Review

The primary aim of this literature study is to identify strategies for insurers and software developers to justify digital price differentiation in terms of fairness and financial returns. We adopt a ‘realist synthesis’ approach, following Turnhout et al. [44]. This approach is characterized by: (1) its focus on addressing the practical problem of how to create solutions for real-world challenges (referred to as the ‘problem in context’) and (2) its emphasis on observable elements, namely interventions (strategies in our terminology) and outcomes (in our case, the justification of digital price differentiation in terms of fairness and financial returns). The remainder of this section details the scoping, query development, literature selection, extraction of strategies, organization of findings, and validation through feedback from practice.

3.2 Scoping the Review

To scope the literature review (and research question) we conducted a mall survey on 36 subjects of professionals in the insurance sector to ensure we focus on the relevant values for digital price differentiation. Appendix B details our findings. First, we found that fairness and wealth are deemed important for justifying digital price differentiation. Second, fairness is deemed slightly but significantly more important. Third, solidarity is relevant for our literature selection and is included in our query. Fourth, transparency and explainability are mentioned as important but considered to be instrumental to fairness. Their importance is reflected in several strategies (e.g., Strategy 25 in Table 4) and further explored in Section 5. Fifth, other mentioned ‘values’ are either encompassed by fairness and wealth or do not refer to actual abstract values (as defined by Oosterlaken [37]). We used these findings to formulate our search query.

3.3 Query Development

To formulate the search query, we identified key terms derived from our research questions. The final query was structured as follows:

```
TITLE-ABS-KEY ( "fair*" OR "solidarity" OR "discriminat*" OR "
    antidiscriminat*" OR "non-discriminat*" OR "equal*" )
AND TITLE-ABS-KEY ( "premium_differentiation" OR ( ( "pric*" AND ( "insur*"
    OR "actuarial" ) ) ) )
```

```

AND TITLE-ABS-KEY ( "algorithm*" OR "machine_learning" OR "AI" OR "
    artificial_intelligence" OR "generalized_linear_model" OR "extreme_
    gradient_boosting" )
AND PUBYEAR > 1999 AND PUBYEAR < 2025
AND ( LIMIT-TO ( DOCTYPE , "ar" ) OR LIMIT-TO ( DOCTYPE , "cp" ) OR LIMIT-
    TO ( DOCTYPE , "ch" ) )
AND ( LIMIT-TO ( LANGUAGE , "English" ) )

```

This query aimed to produce a manageable number of results, avoiding both an insufficient (<5) and an overwhelming (>60) number of papers. Earlier iterations, such as:

```

TITLE-ABS-KEY ( "fair*" ) AND TITLE-ABS-KEY ( "premium_differentiation" )
AND TITLE-ABS-KEY ( "generalized_linear_model" )

```

```

TITLE-ABS-KEY ( "fair*" ) AND TITLE-ABS-KEY ( "premium_differentiation" )

```

yielded too few results. Applying the final query resulted in 54 papers.

Key considerations in the query formulation included:

Use of Synonyms for Fairness Multiple terms were included to ensure comprehensive coverage. ‘Solidarity’ was added based on the results in Section 3.2, while ‘equal*’ and ‘discriminat*’ were incorporated based on domain knowledge and iterative refinements to maintain a pragmatic number of results.

Focus on Insurance The search specifically targeted terms related to insurers and price differentiation to ensure relevance.

Use of Synonyms for digital price differentiation digital price differentiation refers to the use of machine learning algorithms to assess financial risk per client [52]. Since computational methods for price differentiation often overlap with algorithmic and AI-based approaches, these terms were included alongside specific techniques used in practice, such as generalized linear models and extreme gradient boosting.

Omitting Wealth We did not restrict our search to papers explicitly mentioning ‘wealth’ or synonyms like ‘profitability’ or ‘returns.’ As discussed in Section 3.2, optimizing returns is generally assumed in practice. Papers addressing digital price differentiation with a fairness focus were deemed sufficient for covering both aspects.

To ensure quality and consistency, only peer-reviewed, English-language papers were selected. The final query was executed in Scopus on October 18, 2024, yielding 54 papers.

3.4 Selection Process, Exclusion Criteria, and Additional Papers

Figure 1 illustrates the selection process for the literature review. The 54 papers retrieved via Scopus were screened against the following exclusion criteria, leading to the removal of 25 papers:

Different Type of Insurance Papers focusing on niche insurance domains, such as weather insurance for wind farms or agricultural insurance, were excluded.

Not Applied and Only Mathematical Papers that did not apply digital price differentiation in a real-world context but instead focused solely on mathematical proofs were excluded.

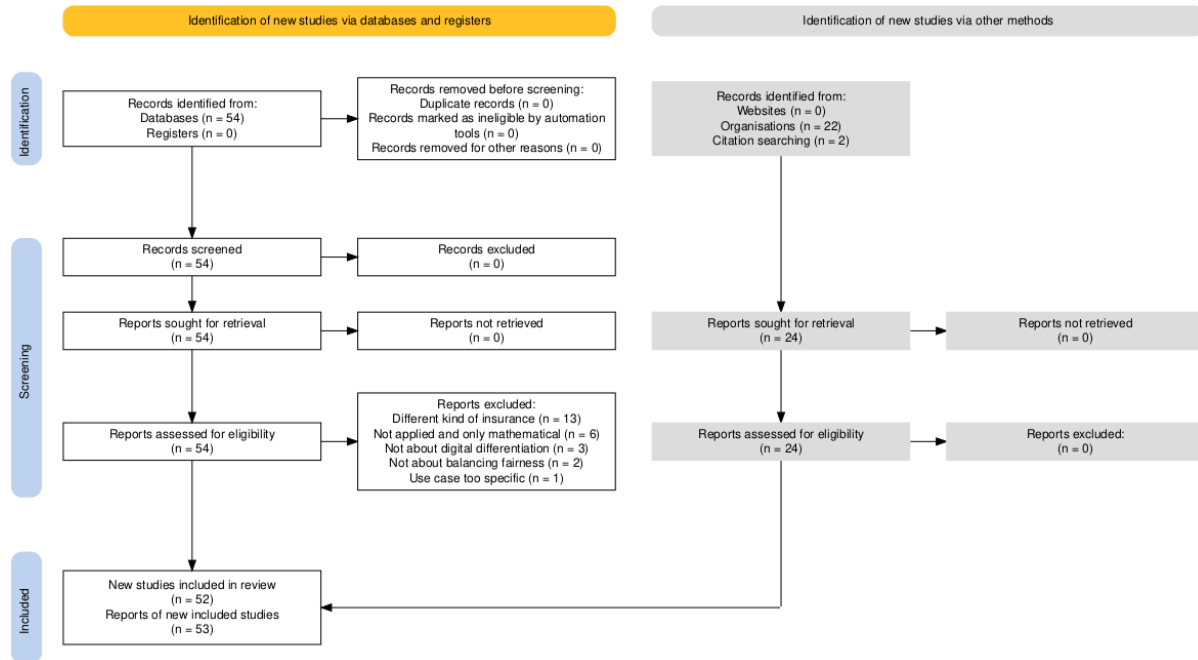


Fig. 1. PRISM selection process to distill papers.

Not About digital price differentiation Papers that discussed price differentiation without involving algorithms or computational models were excluded.

Not About Balancing Fairness Papers that optimized wealth through algorithmic techniques without addressing fairness considerations were excluded.

Use Case Too Specific Papers focusing on very specific regional studies (e.g., a single study on insurance in Thailand) were excluded.

These criteria ensured that our selection remained focused on the socio-technical challenge of how insurers and software developers can justify digital price differentiation. We prioritized papers discussing justification strategies rather than those purely analyzing algorithmic optimization techniques.

In addition to the 28 academic papers selected, 24 additional papers were identified through academic experts and practice partners (n=22) and citation searching (n=2). Practice partners comprise a core group of two software makers for insurance, one insurance company and a Dutch insurance association, and a (peripheral) support group of one additional insurance company, one software maker, one insurance consultancy company and one independent regulator. Academic experts comprise 3 scientist working at the (applied) university on fairness and AI.

At the beginning of our research, we contacted practice partners and academic experts with our main research question and requested relevant literature. Additional papers were suggested from our core group during workshops and discussions related to the broader research project. These non-academic sources, including manuals and protocols, provided valuable practical insights. Combining the academic papers (n=28) and practice-oriented papers

(n=24), we obtained a final set of 52 papers for extracting strategies to justify digital price differentiation in terms of both fairness and financial returns.

3.5 Constructing and Arranging The List

Extracting strategies and synthesise these in the resulting list in Tables 4-6 was an iterative proces reading papers, noting down strategies, ordering these in categories, merging strategies that are similar and redefining and re-ordering the phases and categories. To define the categories [34] was used as a basis. Similarly to Madaio et al. [34], we use temporal phases in our construction and have a determination ('define' in Madaio et al. [34]) and adjusting phase ('build' in Madaio et al. [34] but in our case a model already exists and mainly needs to be adjusted to improve fairness). Unlike Madaio et al. [34], we focus on five phases (1) *Understand* (2) *Determine* (3) *Adjust* (4) *Evaluate* and (5) *Communicate*, making both understanding and communication an explicit phase. Creating understanding within an organization is a necessary prerequisite to enable cooperation on balancing fairness and wealth and focus on evaluating and communicating the model. The evaluation is based on the crisp-dm cycle and the fairness handbook. Communication is a key social factor in justifying the choices mentioned in our literature selection but often overlooked. A main focus of the work was combining and reformulating strategies to make them mutually exclusive and complete in relation to our papers. Our main choices regarding this were on the granularity of presenting the strategies, the wording and separating the legal and ethical aspects (Table 8) and referring to them in other strategies.

We conducted semi-structured interviews with our 4 core practice partners (i.e., one policy adviser from a Dutch insurer association, one actuary and manager from an insurer software company, and one ML expert from the same company, see section 3.5). The interviews, held in Dutch on MS Teams, included a list of strategies with an additional explanation column (see Appendix B.2). For each strategy, we asked: (1) Do you understand the strategy? (2) If not, does the explanation help? (3) If not, why? We also inquired if the list covered all important aspects for justifying digital price differentiation.

Most strategies (22) were understood immediately, with a few minor rewording suggestions. Sixteen were fully understood, and ten required further explanation but were understood after that. Participants suggested shortening the list for ease of use but had no further improvements. We used their feedback to rephrase, remove, or merge strategies and align terminology with practitioner language, such as changing 'subsidiarity' to 'necessity' (see Table 8). Domain knowledge was also integrated, such as emphasizing information asymmetry between insurers (not customers and insurers) in strategy 9. We also simplified the depth of certain strategies for broader understanding.

We shared the list with academic experts for feedback on wording and completeness. Academics were selected by e-mailing our periphal support group and Dutch (applied) university departments related to AI and Fairness. Technical details, like adjusting interactions between sensitive variables (principle 21, Table 5), were added, along with adjustments to emphasize intersectionality and the self-propagating effects of bias (see the full interactive list, Appendix A).

4 Result

Table 3-8 presents our main contribution: a list of strategies for justifying digital price differentiation in terms of fairness and wealth, organized into five phases: (1) *Understand* (Table 3), (2) *Determine* (Table 4), (3) *Adjust*

(Table 5), (4) *Evaluate* (Table 6) and (5) *Communicate* (Table 7). Additionally, eight legal and ethical principles are referenced throughout these phases, influencing the strategies (Table 8). Principles and strategies thus differ in that the principles are relevant in all phases while strategies are split into five phases.

4.1 List of Strategies Organized into Five Phases

The first phase, *Understanding*, is crucial for justifying fair and profitable price differentiation [33]. This phase stresses the importance of recognizing the nuances of discrimination, different forms of bias, and the critical role of sensitive variables in achieving fairness. By understanding these factors, organizations can balance fairness with predictive accuracy, address legal uncertainties, and mitigate privacy concerns. Moreover, market dynamics like information asymmetry and non-risk-based pricing must be considered to maintain competitiveness while adhering to ethical standards. This understanding ensures alignment, creating clarity and focus on how fairness and wealth can be justified, minimizing internal debates.

In the *Determination* phase, various strategies balance fairness and profitability. For example, determining the target group for insurance (strategy 11) directly affects both fairness and profitability. Smaller, homogeneous groups may reduce premium disparities but lead to discrimination risks, as discussed by [25]. The handling of sensitive variables (strategy 12) also plays a critical role, with different levels of rigor affecting fairness outcomes [25]. Depending on this level of rigor one needs to investigate which variables are ethically and legally sensitive and which ones are not (strategy 15); in particular, as understanding why variables are sensitive helps in justifying (the extent of) the use of these variables. Including non-risk factors like discounts (strategy 13) is increasingly regulated in markets like the U.S. and U.K. [29, 51], while setting lower bounds for fairness and expected return (strategy 14) can help assess discrimination risks using principles like the "80% rule" [39]. Next, one should determine which fairness measure one uses; strategy 16-19 describe the core choices such as the choice between individual and group fairness and whether one directly uses a protected variable to measure fairness (or a proxy)[2, 40]. Finally, fairness mitigation strategies (strategy 20) focus on reducing bias at different stages, making determination a pivotal step in justifying fair and profitable price differentiation.

The *Adjustment* phase consists of strategies aimed at mitigating bias and enhancing fairness through pre-processing, in-training, and post-processing techniques. In pre-processing, strategies such as removing sensitive variables and their correlated interactions (strategy 21) or adjusting correlations around sensitive variables (strategy 22) help reduce bias. However, these techniques might reduce fairness and obscure the information of the original data [6, 36, 51]. Another method, adding synthetic data points to balance sensitive relationships (strategy 23), helps mitigate fairness while keeping original data intact [6]. In-training adjustments, such as optimizing an algorithm for both accuracy and fairness (strategy 24), offer a compromise between performance and fairness, though their implementation can be complex [6, 36]. Furthermore, making the algorithm explainable (strategy 25) improves transparency, ensuring risk assessments are understandable [6, 51]. Lastly, post-processing strategies (strategy 26) involve transforming predicted risk assessments to improve fairness, which can avoid recalibrating models and keeps the original data model pipeline and dataset intact. However, this strategy's effectiveness depends on the accuracy of the original model and fairness measures [6, 36, 51]. These adjustment strategies are the bread and butter of actually improving the model such that digital price differentiation becomes fair and remains profitable.

The *Evaluation* phase and *Communication* phase focus on assessing the fairness and profitability of the model and communicating these aspects to stakeholders. Evaluation strategies include comparing the original and

modified models to assess fairness and profitability, ensuring that the adjusted model meets predefined thresholds (strategy 27) [51]. It's also crucial to evaluate the pre-launch model against ethical and legal standards (strategy 28) to ensure compliance with societal and legal expectations [43]. After launch, it is important to monitor the model's fairness through audits and statistical evaluations (strategy 29), verifying that it remains fair and non-discriminatory [6]. On the communication side, it's vital to explain the role of insurance in risk spreading to customers (strategy 30), emphasizing the social aspect [25]. Insurers must also communicate the necessity of risk assessment in determining fair premiums (strategy 31) [25]. Furthermore, fairness decisions should be transparently communicated to regulators, explaining sensitive variables, fairness measures, and model thresholds (strategy 32), ensuring accountability and regulatory compliance [3, 6, 21, 25, 32, 43, 50]. Finally, enabling stakeholders such as customers, competitors, and regulators to challenge the model (strategy 33) promotes transparency and shared responsibility [7].

4.2 The Ethical and Legal Principles

The principles for evaluating the legality and ethics of price differentiation are applied through various strategies across different phases of the pricing process. These principles consist of two parts: (1) a three-step framework to determine the legality of a distinction, and (2) five ethical criteria to assess the sensitivity of a variable. The legal assessment criteria are grounded in regulations concerning indirect discrimination. To determine whether a distinction constitutes indirect discrimination, one must first assess whether the differentiation serves a legitimate aim. In the context of insurance, this usually pertains to the provision of essential services and the accurate assessment of risk (i.e., underwriting) [25, 35, 50] (Principle 1).

Second, the necessity of the variable in achieving this goal is evaluated—often through empirical methods to verify whether it contributes to more accurate risk classification [50]. Necessity is also considered by comparing the variable in question to possible alternatives that might yield comparable results with less negative impact on fairness [21, 35, 50]. The third step involves assessing proportionality, which ensures that the insurer's actions are not excessively burdensome relative to the interests of the insured [35, 50]. This step resembles a cost-benefit analysis, weighing the gains in predictive accuracy against the social or ethical costs of grouping individuals based on specific variables.

While insurers primarily focus on legal compliance, five ethical principles can further illuminate why certain variables are considered sensitive, and to what extent their use may be justifiable. First, the more mutable or influenceable a variable is, the less ethically sensitive it tends to be. For example, gender is generally more ethically sensitive than the number of claim-free years. Second, the stronger the correlation between a variable and the relevant risk—and the weaker its correlation with protected characteristics—the less sensitive it is. Comparing correlations effectively amounts to assessing proportionality (Principle 3).

Third, causal relationships offer a plausible, comprehensible rationale for using a variable. For instance, College voor de Rechten van de Mens [11] argues that a well-established causal link between education level and mortality can justify the use of education as a rating factor in life insurance. Fourth, variables associated with a history of social injustice or discrimination are considered more sensitive (e.g., skin color versus, say, foot size). Fifth, if the use of a variable for differentiation negatively influences individual behavior or broader societal outcomes, it becomes more ethically sensitive. An example is the decision not to use DNA information in life insurance pricing, as this could discourage voluntary DNA registration and hinder scientific progress.

In sum, these five ethical principles help assess the defensibility of using specific variables—or proxies thereof—in pricing decisions.

4.3 Papers not linked to strategies

Tables 3–8 list all papers that propose strategies, as indicated in the 'source' column. However, not every paper identified in the literature review could be directly linked to a specific strategy. Some works focus primarily on technical implementations or broader conceptual issues. The following section describes papers that were reviewed but not explicitly associated with any particular strategy.

4.3.1 *Technical papers focused on algorithmic approaches*

Several papers explore how specific modeling techniques can improve pricing accuracy and, in some cases, fairness. Aruk et al. [5] demonstrate how incorporating location-based risk through Markov-modulated tree-based gradient boosting can lead to fairer pricing outcomes. Kshirsagar [31] compares two modeling approaches and concludes that machine learning techniques outperform traditional Bayesian models, particularly for predicting risk in concession groups. Anzilli [4] shows how the use of fuzzy variables can enhance pricing accuracy. Saleiro et al. [40] introduce Aequitas, an open-source audit toolkit designed to facilitate the evaluation of fairness metrics across different population subgroups. Pe na-Sanchez [38] focuses on data-scarce insurance markets, such as in the Philippines, and demonstrates how the application of a generalized linear model with a Tweedie compound Poisson–Gamma distribution can improve pricing for bundled micro-insurance products.

4.3.2 *Papers on practical applications in insurance pricing*

Cunha and Bravo [13] investigate the value of telematics data in enhancing risk estimation and, consequently, pricing. Their work emphasizes the technical potential for price discrimination, though it does not directly address fairness concerns. Fabris et al. [22] conduct an audit of pricing algorithms used in the Italian car insurance industry. Based on an analysis of 2,160 driver profiles, they find that sensitive attributes such as birthdate and gender are used directly in pricing algorithms, a practice that violates existing regulations.

4.3.3 *Alternative perspectives on fairness*

Some studies approach fairness from non-outcome-based perspectives. Grgič-Hlača et al. [27] emphasize procedural fairness. They assess the acceptability of using certain features based on survey responses, showing how public perception can guide the fair use of variables. Donahue and Barocas [16] use game theory to theoretically examine the trade-off between solidarity and actuarial fairness—highlighting the tension between socially equitable pricing and precision in risk assessment.

5 Discussion, Conclusion & Future Work

This paper has addressed the challenge of ensuring fairness while maintaining financial returns in digital price differentiation within the insurance industry. The main contribution — summarised in Tables 3–8 — is a categorised and ordered set of 33 strategies (and 8 ethical-legal principles employed in several of these strategies). These strategies offer a practical guide for insurers and software developers to navigate the socio-technical dilemma of balancing fairness and profitability. This list is scientifically distinctive due to its comprehensive scope, its

integration of both technical and social dimensions, and its focus on digital price differentiation as a regression problem rather than a classification problem. Thus, our work contributes to the literature on AI and fairness by moving beyond abstract fairness metrics and offering actionable strategies that align with organisational processes and values.

From a practical perspective, professionals in the insurance sector—including data scientists, actuaries, auditors, compliance officers, policymakers, and communications teams—can use this strategy list in three ways: (1) as a guide for selecting relevant strategies, (2) as a checklist to ensure that all necessary steps have been taken, and (3) as a communication tool to structure internal discussions about fairness and financial objectives. For example, the list provides a comprehensive overview of strategies, while the online tool and resources (see Appendix A) help users identify and implement those most applicable to their context. In short, the strategy list supports practice as a reference guide, checklist, and facilitation tool for cross-functional dialogue.

This review has centered on two human values: fairness and wealth. In line with value-sensitive design principles [26], we approached these values not as opposing forces but as dimensions that can be reconciled through thoughtful strategy selection. Although the survey (see Appendix B) could be improved through refinements such as a 7-point Likert scale, a broader and more representative sample, and additional open-ended questions, it confirmed the central importance of both values in digital price differentiation.

The value of wealth is especially interesting due to its multiple interpretations. On one hand, it refers to the financial stability of insurance companies—something required by regulators and essential for the industry’s long-term viability and the public interest. On the other hand, wealth can be philosophically contested as a human value. While not a terminal value in Schwartz’s theory, it is part of the broader value of power [41] and is commonly used in practice to contrast with fairness. Since our aim is to remain close to industry language while acknowledging philosophical nuance, we adopt wealth as a pragmatic proxy for financial returns. Similarly, although terms like solidarity appeared in our literature search, we have treated solidarity as instrumental to fairness in this paper. Future research should explore how both wealth and solidarity are understood and valued in the insurance context.

By combining a systematic literature review with empirical insights from industry professionals, we have tailored this strategy set to the specific needs of the insurance sector, ensuring both relevance and applicability. The findings offer valuable guidance for policymakers, auditors, actuaries, and developers working with algorithmic systems in insurance.

Future research should focus on validating and refining this strategy framework with a larger and more diverse group of professionals. Additional studies could investigate how these strategies can be operationalized in practice and adapted for use in other domains where AI is used for pricing and decision-making.

References

- [1] Andrew Abbott. 1983. Professional Ethics. , 855–885 pages.
- [2] Gemeente Amsterdam. 2022. *The Fairness Handbook*. Technical Report May. Gemeente Amsterdam. 68 pages. <https://drive.google.com/file/d/18FTdKvvvmnmpAwagJIg6YwAa6KEU-Km/view>
- [3] Katrien Antonio. 2022. *Bias, fairness and discrimination-free insurance pricing*. Technical Report. (slides) pages.
- [4] Luca Anzilli. 2012. A possibilistic approach to evaluating equity-linked life insurance policies. *Communications in Computer and Information Science* 300 CCIS, PART 4 (2012), 44 – 53. https://doi.org/10.1007/978-3-642-31724-8_6
- [5] Dennis Arku, Kwabena Doku-Amponsah, and Nathaniel K Howard. 2020. A Markov-modulated tree-based gradient boosting model for auto-insurance risk premium pricing. *Risk and Decision Analysis* 8, 1-2 (2020), 1 – 13. <https://doi.org/10.3233/RDA-180050>

- [6] Solon Barocas, Moritz Hardt, and Arvind Narayanan. 2020. Fairness and Machine Learning. 2019 (2020), 9–35. <https://fairmlbook.org/%0Ahttps://fairmlbook.org>
- [7] Laurence Barry and Arthur Charpentier. 2023. Melting contestation: insurance fairness and machine learning. *Ethics and Information Technology* 25, 4 (2023). <https://doi.org/10.1007/s10676-023-09720-y>
- [8] Emily M Bender and Batya Friedman. 2018. Data Statements for Natural Language Processing: Toward Mitigating System Bias and Enabling Better Science. , 587–604 pages.
- [9] CBS. 2024. ICT-gebruik bij bedrijven; bedrijfstak, 2024.
- [10] Alberto Cevolini and Elena Esposito. 2020. From pool to profile: Social consequences of algorithmic prediction in insurance. *Big Data and Society* 7, 2 (2020). <https://doi.org/10.1177/2053951720939228>
- [11] College voor de Rechten van de Mens. 2014. Advies aan Dazure B.V. over premiedifferentiatie op basis van postcode bij de Finvita overlijdensrisicoverzekering | Mensenrechten. (2014). <https://mensenrechten.nl/nl/publicatie/19173>
- [12] Henriette Cramer, Jean Garcia-Gathright, Sravana Reddy, Aaron Springer, and Romain Takeo Bouyer. 2019. Translation, Tracks & Data: an Algorithmic Bias Effort in Practice. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. ACM, 1–8. <https://doi.org/10.1145/3290607.3299057>
- [13] Lourenço Cunha and Jorge M Bravo. 2022. Automobile Usage-Based-Insurance: Improving Risk Management using Telematics Data. In *Iberian Conference on Information Systems and Technologies, CISTI*, Vol. 2022-June. <https://doi.org/10.23919/CISTI54924.2022.9820146>
- [14] Luca Deck, Jan-Laurin Müller, Conradin Braun, Dominique Zipperling, and Niklas Kühl. 2024. Implications of the AI Act for Non-Discrimination Law and Algorithmic Fairness. *EWAf'24: European Workshop on Algorithmic Fairness* (2024). <http://arxiv.org/abs/2403.20089>
- [15] Luca Deck, Jan-Laurin Müller, Conradin Braun, Dominique Zipperling, and Niklas Kühl. 2024. Implications of the AI Act for Non-Discrimination Law and Algorithmic Fairness. *EWAf'24: European Workshop on Algorithmic Fairness* (2024). <http://arxiv.org/abs/2403.20089>
- [16] Kate Donahue and Solon Barocas. 2021. Better together?: How externalities of size complicate notions of solidarity and actuarial fairness. In *FAccT 2021 - Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*. 185 – 195. <https://doi.org/10.1145/3442188.3445882>
- [17] DrivenData. 2019. Deon: An ethics checklist for data scientists. <http://deon.drivendata.org/>
- [18] EDPB. 2024. Guidelines 1 / 2024 on processing of personal data based on Adopted on 8 October 2024. , 37 pages.
- [19] EIOPA. 2025. *On Opinion on Artificial Intelligence Governance and Risk*. Technical Report February. 22 pages.
- [20] European EU Law. 2004. Council Directive 2004/113/EC of 13 December 2004 implementing the principle of equal treatment between men and women in the access to and supply of goods and services.
- [21] European EU Law. 2016. Consolidated text: Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 9. <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX%3A02016R0679-20160504>
- [22] Alessandro Fabris, Alan Mishler, Stefano Gottardi, Mattia Carletti, Matteo Daicampi, Gian Antonio Susto, and Gianmaria Silvello. 2021. Algorithmic Audit of Italian Car Insurance: Evidence of Unfairness in Access and Pricing. In *AIES 2021 - Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*. 458 – 468. <https://doi.org/10.1145/3461702.3462569>
- [23] Mary Flanagan and Helen Nissenbaum. 2014. Values at Play in Digital Games.
- [24] Johns Hopkins Center for Government Excellence. 2019. Ethics & Algorithms Toolkit. <http://ethicstoolkit.ai/>
- [25] Edward W. Frees and Fei Huang. 2023. The Discriminating (Pricing) Actuary. *North American Actuarial Journal* 27, 1 (2023), 2–24. <https://doi.org/10.1080/10920277.2021.1951296>
- [26] Batya Friedman. 1996. Value-sensitive design. *Interactions* 3, 6 (1996), 16–23. <https://doi.org/10.1145/242485.242493>
- [27] Nina Grgić-Hlača, Muhammad Bilal Zafar, Krishna P Gummadi, and Adrian Weller. 2018. Beyond distributive fairness in algorithmic decision making: Feature selection for procedurally fair learning. *32nd AAAI Conference on Artificial Intelligence, AAAI 2018* (2018), 51–60. <https://doi.org/10.1609/aaai.v32i1.11296>
- [28] European Union High-level Expert Group. 2019. Ethics Guidelines for Trustworthy AI: Building trust in human-centric AI. <https://ec.europa.eu/futurium/en/ai-allianceconsultation/guidelines>

- [29] R. Guy Thomas. 2012. Non-risk price discrimination in insurance: Market outcomes and public policy. *Geneva Papers on Risk and Insurance: Issues and Practice* 37, 1 (2012), 27–46. <https://doi.org/10.1057/gpp.2011.32>
- [30] Dutch Association of Insurers. 2018. *Code of Conduct for*. Technical Report.
- [31] Rohun Kshirsagar, Li-Yen Hsu, Charles H Greenberg, Matthew McClelland, Anushadevi Mohan, Wideet Shende, Nicolas P Tilmans, Min Guo, Ankit Chheda, Meredith Trotter, Shonket Ray, and Miguel Alvarado. 2021. Accurate and Interpretable Machine Learning for Transparent Pricing of Health Insurance Plans. In *35th AAAI Conference on Artificial Intelligence, AAAI 2021*, Vol. 17A. 15127 – 15136. <https://www.scopus.com/inward/record.uri?eid=2-s2.0-85130091421&partnerID=40&md5=51f5ae14e156ff4a9a651956da8c8518>
- [32] M. Lindholm, R. Richman, A. Tsanakas, and M. V. Wüthrich. 2022. Discrimination-Free insurance pricing. *ASTIN Bulletin* 52, 1 (2022), 55–89. <https://doi.org/10.1017/asb.2021.23>
- [33] Michael A Madaio, Jingya Chen, Hanna Wallach, and Jennifer Wortman Vaughan. 2024. Tinker, Tailor, Configure, Customize: The Articulation Work of Contextualizing an \{AI\} Fairness Checklist. *Proc. ACM Hum.-Comput. Interact.* 8, CSCW1 (2024), Article 214.
- [34] Michael A Madaio, Luke Stark, Jennifer Wortman Vaughan, and Hanna Wallach. 2020. Co-Designing Checklists to Understand Organizational Challenges and Opportunities around Fairness in AI. *Conference on Human Factors in Computing Systems - Proceedings* (2020), 1–14. <https://doi.org/10.1145/3313831.3376445>
- [35] Ministerie van Binnenlandse Zaken en Koninkrijksrelaties. 2020. Algemene Wet Gelijke Behandeling. <https://wetten.overheid.nl/jci1.3:c:BWBR0006502&z=2020-01-01&g=2020-01-01>
- [36] Mulah Moriah, Franck Vermet, and Arthur Charpentier. 2024. Measuring and mitigating biases in motor insurance pricing. *European Actuarial Journal* (2024). <https://doi.org/10.1007/s13385-024-00390-8>
- [37] Ilse Oosterlaken. 2015. *Handbook of Ethics, Values, and Technological Design*. 221–250 pages. <https://doi.org/10.1007/978-94-007-6970-0>
- [38] Inmaculada Peña-Sanchez. 2019. Applying the Tweedie model for improved microinsurance pricing. *Geneva Papers on Risk and Insurance: Issues and Practice* 44, 3 (2019), 365 – 381. <https://doi.org/10.1057/s41288-019-00130-0>
- [39] Ronald B Rubin. 1979. The Uniform Guidelines on Employee Selection Procedures : Compromises and Controversies.
- [40] Pedro Saleiro, Benedict Kuester, Loren Hinkson, Jesse London, Abby Stevens, Ari Anisfeld, Kit T. Rodolfa, and Rayid Ghani. 2018. Aequitas: A Bias and Fairness Audit Toolkit. 2018 (2018). <http://arxiv.org/abs/1811.05577>
- [41] Shalom H Schwartz. 2012. An Overview of the Schwartz Theory of Basic Values. *Online Readings in Psychology and Culture* 2 (2012), 1–20. <https://doi.org/http://dx.doi.org/10.9707/2307-0919.1116>
- [42] George A Steiner. 2010. *Strategic planning*. Simon and Schuster.
- [43] Jelte Timmer, Isabel Elias, Linda Kool, and Rinie van Est. 2015. *Berekende risico's: Verzekeren in de datagedreven samenleving*.
- [44] Koen van Turnhout, D G Andriessen, and Petra Cremers. 2023. *Handboek ontwerpgericht wetenschappelijk onderzoek : ontwerpend onderzoeken in sociale contexten*. Boom uitgevers, Amsterdam. 392 pages. <https://doi.org/LK-https://hu.on.worldcat.org/oclc/1378467861>
- [45] Utrecht Data School. 2021. Impact Assessment Mensenrechten en Algoritmes. , 95 pages.
- [46] Utrecht Data School. 2022. De Ethische Data Assistent.
- [47] Marvin van Bekkum and Frederik Zuiderveen Borgesius. 2023. Using sensitive data to prevent discrimination by artificial intelligence: Does the GDPR need a new exception? *Computer Law and Security Review* 48 (2023). <https://doi.org/10.1016/j.clsr.2022.105770>
- [48] Verbond voor Verzekeraars. 2023. *Solidariteitsmonitor*. Technical Report. 59 pages.
- [49] Sahil Verma and Julia Rubin. 2018. Fairness definitions explained. In *Proceedings of the International Workshop on Software Fairness*. ACM, New York, NY, USA, 1–7. <https://doi.org/10.1145/3194770.3194776>
- [50] College voor de Rechten van de Mens. 2025. Toetsingskader risicoprofilering. (2025).
- [51] Xi Xin and Fei Huang. 2024. Antidiscrimination Insurance Pricing: Regulations, Fairness Criteria, and Models. *North American Actuarial Journal* 28, 2 (2024), 285–319. <https://doi.org/10.1080/10920277.2023.2190528>
- [52] Frederik J. Zuiderveen Borges. 2022. Digitale discriminatie en differentiatie : het recht is er nog niet klaar voor.
- [53] Salih Tayfun Ince. 2022. European Union Law and Mitigation of Artificial Intelligence-Related Discrimination Risks in the Private Sector: With Special Focus on the Proposed Artificial Intelligence Act. *Annales de la Faculte de Droit d'Istanbul* 71 (2022), 265 – 307. <https://doi.org/10.26650/annales.2022.71.0002>