



Safe and Effective AI in Regulated Debt Advice

Data Preparation Guide

Exploring the role of generative AI-enabled self-serve tools in a complex, vulnerability-rich and regulated environment

This guide helps organisations that hold substantial data assets understand how to prepare, govern and use their data responsibly for AI applications in regulated debt advice. It is a companion to the evidence report and safe-use framework.

Data Preparation Guide v1.0
2026

Contents

Contents

1	Why this guide?	3
2	Identifying your data assets	4
3	Assessing Data Quality	7
4	Preparing Your Data for Each Level.....	9
5	Data Governance for AI	14
6	Technical Considerations.....	16
7	Getting Started: A Practical Roadmap.....	18

1 Why this guide?

Large debt advice organisations hold some of the richest data on personal financial difficulty in the UK. Case records spanning years or decades, financial statements across hundreds of thousands of clients, creditor correspondence, quality assurance findings, outcome data and specialist knowledge bases represent a significant asset.

As generative AI develops, these data assets become strategically significant. An AI tool trained or configured using an organisation's own data could, in principle, provide guidance that reflects the organisation's specific approach, standards and client population. However, the path from raw organisational data to a safe, effective AI tool is not straightforward. It requires careful assessment of data quality, rigorous governance, appropriate anonymisation, technical preparation and ongoing maintenance.

This guide addresses the practical questions that organisations face when considering this path: what data do you have, what is it suitable for, how do you prepare it, what are the risks and what governance is needed. It is grounded in findings from sector research conducted for this project, which identified both the potential and the significant challenges of using organisational data for AI in regulated debt advice. The guide is structured around the Debt Advice AI Roadmap, which provides the framework for understanding which applications your data might support.

2 Identifying your data assets

Before considering any AI applications, organisations need a clear picture of what data they hold, where it sits and what condition it is in. The research identified the following categories of data commonly held by debt advice organisations, along with their potential applications and risks.

Case management records

Most organisations maintain structured case management systems. The research identified several in use across the sector, including Casebook (widely used by Citizens Advice), AdvicePro, Petra and bespoke CRM systems developed by individual organisations. These systems typically contain client demographic information, debt details, income and expenditure data, actions taken, outcomes recorded and follow-up notes.

Case management data is the single richest potential source for AI training, but also the most sensitive. It contains identifiable personal information, including financial, health and vulnerability disclosures. The volume varies enormously: some organisations hold tens of thousands of case records; larger national providers hold hundreds of thousands or more.

Financial statements

Income and expenditure statements, commonly using the Standard Financial Statement (SFS) or Common Financial Statement (CFS) formats, represent a particularly valuable data category. These contain detailed household budget information, income from all sources and expenditure across standardised categories.

The research revealed significant complexity in the construction of financial statements. Advisers described the process as an art form requiring professional challenge, with figures negotiated rather than recorded. This means that financial statement data reflects both the client's actual circumstances and the adviser's professional interpretation. AI tools trained on this data will learn both dimensions.

Caution: Financial statements may contain figures that were adjusted for tactical purposes (e.g. to support a particular debt solution application). Training AI on these statements without understanding this context could produce outputs that reproduce tactical adjustments as standard practice.

Full data source inventory

The following table maps the data sources identified through the research to their potential AI applications and primary risks:

Data source	What it contains	Applicable levels	Key risks
Case management records	Client demographics, debt details, income/expenditure, actions, outcomes, follow-up notes	1.1, 1.2, 1.3, 2.3, 2.4	Highly sensitive Personally Identifiable Information (PII). Consent limitations. Historical bias risk
Financial statements	Detailed household budgets, income sources, expenditure categories, surplus/deficit	1.2, 1.3, 2.4	Contains professional interpretation alongside facts. Tactical adjustments may be present
Case notes and narratives	Adviser-written records of conversations, assessments, reasoning, vulnerability observations	1.1, 1.2, 1.3	Unstructured text. May contain sensitive disclosures. Quality varies by adviser
Call recordings and transcripts	Audio recordings and AI-generated transcripts of adviser-client conversations	1.1, 1.3	Most detailed record of practice. Consent required. Storage-intensive. Transcription errors
Creditor correspondence	Letters, emails and notes from creditor interactions, including offers, rejections, and arrangements	1.1, 2.5	Contains creditor-specific intelligence. May be outdated rapidly
Knowledge bases and guidance	Internal adviser guidance, policy documents, procedure manuals, legislation summaries	1.2, 2.1, 2.2	Highest value for 2.1/2.2 tools. Must be current. Jurisdictional accuracy critical
Quality assurance records	QA assessments, compliance scorecards, supervisor feedback, error patterns	1.2, 1.3	Defines what good practice looks like. May reflect changing standards over time
Template libraries	Letter templates, standard paragraphs, creditor-specific wording, procedural checklists	1.1, 2.5	Efficient starting point. Must be current. Jurisdictional variants needed
Open banking data	Bank transaction data accessed with client consent, categorised spending, and income verification	1.3, 2.4	Consent-specific. Categorisation is often inaccurate. Does not capture cash transactions
Credit reference data	Experian/Equifax reports accessed during the advice process showing credit commitments, defaults	2.3, 2.4	Third-party data. Licensing restrictions. Point-in-time accuracy only
Outcome and impact data	Solution types implemented, completion rates, re-contact rates, and client satisfaction	1.3	Measures what happened, not why. May not capture long-term outcomes
Funder reporting data	MaPS outcomes framework data, KPI reports, demographic breakdowns	1.3	Structured, but may not capture advice quality. Reporting definitions change over time
Training materials	Adviser training content, case studies, assessment	1.2	Valuable for knowledge support tools. Must be updated as practice evolves

	scenarios, and competency frameworks		
Website and digital content	Published guidance, self-help resources, FAQ content, chatbot logs, if any	2.1, 2.2	Already public. Good starting point for 2.1 tools. Must be current

3 Assessing Data Quality

The value of any AI tool is fundamentally constrained by the quality of data it draws from. Before investing in AI development, organisations should conduct an honest assessment of their data quality across several dimensions.

Completeness

How complete are your case records? The research found that case note quality varies significantly between advisers, offices and time periods. Some advisers produce detailed, structured notes; others produce minimal records. Cases that were resolved quickly may have sparse documentation, while complex cases may be thoroughly recorded. AI tools trained on incomplete data may produce outputs that reflect the gaps rather than the full picture.

Accuracy

How accurate is the information recorded? Financial statement data is subject to the professional interpretation described throughout this research. Advisers challenge self-reported income and expenditure figures, but they may still contain inaccuracies. Creditor information changes over time. Outcome data may not reflect long-term results. Any AI application must account for these inherent accuracy limitations.

Timeliness

How current is your data? Debt advice operates in a fast-changing regulatory environment. Case records from three years ago may reference debt solutions, fee structures, eligibility criteria or enforcement practices that have since changed. Knowledge bases and guidance documents require continuous updating. The research found that advisers spend a lot of time remaining current in fast-changing areas.

Consistency

How consistently is data recorded across your organisation? The research revealed significant variation in how different offices, teams and individual advisers record information. Different advisers may use different terminology, levels of detail, and approaches to documenting the same type of case. AI tools trained on inconsistent data may produce variable-quality outputs.

Representativeness

Does your data represent the full range of clients and circumstances you serve? If certain client groups are underrepresented in your data (for example, clients with specific vulnerabilities, clients from particular demographic backgrounds or clients with certain debt types), AI tools trained on that data may perform poorly for those groups. This is a significant equity concern.

Consent and legal basis

Under what legal basis was the data collected, and does that basis extend to AI training? Most client data in debt advice was collected to provide advice. Using it to train AI systems may require a different legal basis under UK GDPR. Organisations must assess whether existing consent covers AI use and, where it does not, determine how to proceed lawfully.

4 Preparing Your Data for Each Level

Different AI applications require different data preparation approaches. This chapter maps the Debt Advice AI Roadmap to specific data preparation steps, organised by the three levels. As with the Roadmap itself, not every organisation will need to prepare data for all three levels. Start with your identified use case and prepare only the data that supports it.

Level 1: Adviser Assistance

Level 1 tools support adviser workflows and organisational operations. Data preparation requirements are the least onerous because the output is consumed by professionals who can assess its quality.

Administrative Assistance (transcription, case notes, letters) (1.1)

Data needed

- Call recordings and existing transcripts for training transcription accuracy
- Case note examples showing your organisation's preferred format, level of detail and terminology
- Letter templates and examples of well-drafted creditor correspondence

Preparation steps

- Compile a representative sample of well-written case notes (minimum 500) that reflect your quality standards
- Assemble your current letter template library with jurisdiction-specific variants
- If using transcription, test against recordings that include specialist terminology, creditor names and accents common in your client base
- Remove or redact client-identifiable information from training examples (names, addresses, account numbers, dates of birth)
- Review sample outputs against your quality standards before full deployment

Knowledge and Compliance Support (1.2)

Data needed

- Internal knowledge bases, policy documents, adviser guidance and procedure manuals
- Legislation summaries and regulatory guidance in your own interpretation
- QA standards and compliance frameworks

Preparation steps

- Audit your knowledge base for currency: identify documents that are out of date or contain superseded information
- Establish version control: every document should have a last-reviewed date and a next-review date
- Tag content by jurisdiction (England/Wales, Scotland or UK-wide)

- Identify content that reflects your organisation's interpretation rather than settled law and flag it accordingly
- Create a maintenance schedule
- Define update triggers: what events (new case law, legislative change, regulatory guidance) require immediate knowledge base updates

Operational Intelligence (1.3)

Data needed

- Anonymised or aggregated case outcome data across the organisation
- Caseload and workload data by team, office and time period
- QA assessment results and trend data
- Client demographic and referral source data

Preparation steps

- Ensure all data used for operational analysis is properly anonymised or aggregated
- Establish baseline metrics before introducing AI analysis to enable a before-and-after comparison
- Check for systematic gaps: if certain offices or teams have lower recording rates, analysis will underrepresent their caseloads
- Define the questions you want operational intelligence to answer before preparing data - data preparation should be purpose-driven

Level 2: Client Self-Serve: Information and Assistance

Level 2 tools interact directly with clients or support them alongside the advice journey. Data preparation requirements increase significantly because outputs affect individuals in potentially vulnerable circumstances. Requirements are cumulative: higher-numbered application types inherit the data standards of those below.

General Information and Contextual Guidance (2.1, 2.2)

Data needed

- Published guidance, website content, FAQ databases and self-help resources
- Internal knowledge bases (carefully reviewed for client-facing suitability)
- Approved terminology glossaries and plain-language explanations

Preparation steps

- Review all content for reading age accessibility
- Remove or rewrite professional jargon, internal terminology and adviser-facing language
- Tag every piece of content by jurisdiction and verify accuracy
- Establish a review cycle: content used by client-facing tools must be checked more frequently than internal knowledge bases (recommended: monthly for 2.1, fortnightly for 2.2)
- Test with sample queries covering the most common client questions identified in your case data
- Establish accuracy benchmarks

Structured Triage (2.3)

Data needed

- Historical triage decision data: how cases were routed and whether routing proved appropriate
- Vulnerability indicator data: which indicators were present and how they were identified
- Urgency assessment data: which cases required emergency response, and what triggered that assessment

Preparation steps

- Analyse your historical triage data to identify patterns in correct and incorrect routing
- Document the implicit decision rules your experienced advisers use for triage (these are rarely written down)
- Map your vulnerability indicators against the escalation trigger list in the Safe-Use Framework
- Test triage logic against a sample of real (anonymised) cases and compare AI routing with human routing
- Calculate the false negative rate: how often does the system fail to escalate cases that should have been escalated? Target: as close to zero as possible

Pre-appointment Data Gathering (2.4)

Data needed

- Financial statement data: completed SFS/CFS forms showing the range of client circumstances
- Common data quality issues: which fields are most often incomplete, inaccurate or misunderstood by clients
- The typical gap between client-submitted data and adviser-validated data

Preparation steps

- Analyse your financial statement data to identify the most common discrepancies between initial submissions and adviser-validated versions
- Build validation rules based on your advisers' challenge patterns (e.g. food expenditure below realistic thresholds, missing irregular expenses, suppressed essential spending)
- Design the tool to flag potential issues for adviser review rather than correcting them automatically
- Ensure the tool explicitly states that an adviser will review the financial statement and may change
- Test with real (anonymised) client data and compare AI-gathered information with adviser-gathered information

Critical warning: Pre-appointment data-gathering tools fall within the core territory of regulated advice. Disposable income calculations directly affect solution eligibility. If the tool produces an unrealistic financial picture, it may set expectations that the adviser must then undo. All financial data gathered by AI must be treated as a starting point for professional review, never as a finished assessment.

Communication Support, Documentation Support and Ongoing Engagement (2.5 - 2.7)

Data needed

- Creditor correspondence examples covering the full range of creditor types and situations
- Sample appointment preparation checklists and document lists
- Engagement and re-engagement communications that have proven effective

Preparation steps

- Compile creditor-specific letter examples showing what works with different organisations
- Remove high-stakes correspondence from the training set (statute-barred debt responses, court submissions, insolvency applications) - AI must not generate these without professional review
- Review all client-facing content for tone: it should be warm, supportive and non-judgemental
- Test re-engagement messages with sensitivity - clients who have disengaged may be in crisis

Level 3: Agentic AI: Taking Action and Sustained Relationships

Level 3 applications take actions or maintain sustained evolving relationships with clients. They require all of the data preparation for the underlying application type, plus additional governance. This level also includes personalised advisory (3.1), which requires agentic capabilities to sustain its context-maintaining relationship model.

Personalised Advisory (3.1)

Not recommended for development without extensive piloting and regulatory engagement. Data preparation for 3.1 would require all Level 2 data preparation plus longitudinal outcome data demonstrating that AI-supported processes produce comparable or better outcomes than fully human-delivered advice. The sector survey finding that only 18% now consider this suitable (12% for all debts) reflects genuine, evidence-based caution.

Adviser-Directed Actions, Semi-Autonomous Actions and Autonomous Process Management (3.2 – 3.4)

Level 3 applications that take actions require all of the data preparation for the underlying application type, plus additional governance:

- Map every possible action the system could take and classify it as routine or consequential
- For each consequential action, define the data the system needs to confirm before proceeding
- Establish audit trail requirements: what data must be logged when an action is taken
- Define rollback data requirements: what information is needed to reverse an action if an error occurs

5 Data Governance for AI

Using organisational data for AI is not simply a technical exercise. It requires a governance framework that addresses legal, ethical and operational dimensions.

Legal requirements

- Complete a Data Protection Impact Assessment (DPIA) before any data is used for AI training or configuration
- Identify the lawful basis for processing under UK GDPR for each data source and each intended use
- Assess whether existing client consent covers AI training (in many cases, it will not)
- Review data sharing agreements with third parties (e.g. open banking providers, credit reference agencies) to confirm AI use is permitted
- Consider whether model outputs could indirectly reveal client-identifiable information even if the training data was anonymised
- Establish data retention policies specific to AI training data

Ethical considerations

The research identified several ethical concerns that organisations must address:

Algorithmic bias: Historical case data may reflect existing inequities. If certain demographic groups have historically received different treatment, AI trained on that data may reproduce those patterns. Organisations should test for differential outcomes across demographic groups before deployment.

Consent and expectations: Clients shared their information to receive advice. Using it to train AI systems changes the nature of the data relationship. Even where legal bases exist, organisations should consider whether this use aligns with reasonable client expectations.

Commercial use: If organisational data is used to train AI tools that are subsequently commercialised or licensed to third parties, this raises additional questions about consent and governance. Several sector survey respondents expressed concern about this scenario.

Data from vulnerable individuals: Case records may contain particularly sensitive disclosures. The inclusion of this data in AI training sets requires especially careful consideration of consent, anonymisation and potential for re-identification.

Operational governance

- Appoint a named data governance lead for AI initiatives
- Establish a data governance board or committee with representation from advice delivery, compliance, data protection and technology
- Create a data catalogue documenting what data exists, where it is stored, who owns it, what condition it is in and what it may be used for
- Define data quality standards for AI use (these may be higher than standards for other purposes)
- Establish monitoring processes: regularly assess whether AI tools built on your data are producing outputs that meet quality, accuracy and fairness standards

- Plan for data currency: establish processes for updating AI training data or knowledge bases as the underlying information changes
- Document all decisions about data use, including what data was included, what was excluded and why

6 Technical Considerations

This chapter provides practical guidance on the technical aspects of preparing data for AI. It is written for non-technical leaders; organisations undertaking AI development will need specialist technical support.

Approaches to using your data

There are three main approaches, each with different data requirements:

Retrieval-Augmented Generation (RAG) is the most common and lowest-risk approach. The AI model retrieves information from your approved knowledge base to inform its responses. Your data is not used to train the model itself; it is used as a reference source. This is the recommended approach for Level 1 and application types 2.1–2.2. It requires well-organised, current and clearly tagged knowledge content, but does not require the large datasets needed for fine-tuning.

Fine-tuning involves further training an existing AI model on your specific data so that it learns your organisation's patterns, terminology and approaches. This produces more tailored outputs but requires larger datasets (typically thousands to tens of thousands of examples) and carries a greater risk of embedding biases or errors present in the training data. It may be appropriate for organisations with substantial, high-quality data pursuing applications at 2.3 or above.

Custom training involves building a model from scratch using your data. This is the most resource-intensive approach and is unlikely to be appropriate for most advice organisations. It requires very large datasets, significant technical expertise and substantial computing resources.

Anonymisation and de-identification

Before using case data for any AI purpose, personal identifiers must be removed or replaced. This includes names, addresses, dates of birth, national insurance numbers, account numbers, telephone numbers, email addresses and any other information that could identify an individual directly or in combination with other data. Even anonymised data can sometimes be re-identified through combinations of characteristics; organisations should assess this risk as part of their DPIA.

Data formatting and structure

- Structure data consistently: ensure case records, financial statements and knowledge bases use consistent formatting
- Use clear metadata: tag all content with creation date, last-reviewed date, jurisdiction, document type and author where applicable
- Handle unstructured data carefully: case notes and narratives contain rich information but require natural language processing to extract structured insights
- Maintain provenance: for every piece of training data, record where it came from, when it was created and under what conditions

Infrastructure and security

- Establish whether AI training and inference will occur on-premises, in a private cloud or via a third-party provider

- Ensure that data in transit and at rest is encrypted to the same standards as other sensitive data
- Confirm that AI infrastructure meets the organisation's information security requirements
- Consider data residency: UK GDPR imposes restrictions on transferring personal data outside the UK

7 Getting Started: A Practical Roadmap

For organisations ready to explore using their data for AI, the following roadmap provides a staged approach aligned with the Debt Advice AI Roadmap. As with the Roadmap itself, progression should be evidence-based and driven by demonstrated success rather than external pressure. Not every organisation will need or want to progress beyond Phase 1.

Phase 1: Foundation

- Conduct a data audit: catalogue what you hold, where it sits and what condition it is in
- Complete a DPIA for the intended AI use
- Establish data governance arrangements (lead, board/committee, policies)
- Assess legal bases for AI use of each data source
- Define your starting level and application type (recommended: 1.1 administrative assistance)
- Identify a pilot use case with clear success criteria

Phase 2: Pilot

- Prepare data for the pilot use case following the level-specific guidance above
- Select or develop the AI tool, working with technology partners who understand the sector
- Test with a controlled group of advisers or users
- Monitor accuracy, error patterns and user experience
- Gather feedback and iterate
- Document lessons learned

Phase 3: Evaluate and Expand

- Assess pilot outcomes against success criteria
- Review data quality issues that emerged during the pilot
- Update data governance processes based on operational experience
- If successful, plan expansion to additional use cases within the same level or cautious progression to adjacent levels
- Establish ongoing monitoring and review cycles

Phase 4: Maturity

- Embed AI data governance into organisational standard practice
- Develop internal capacity for data preparation and quality assurance
- Consider sharing anonymised learnings with the sector to support collective progress
- Regularly reassess data quality, currency and representativeness
- Stay engaged with regulatory developments and sector standards

The golden rule

Data-driven AI development in debt advice requires the same rigour, governance and client-centred thinking that characterises good advice practice itself. If you would not give a piece of information to a client without checking it first, do not let an AI tool do so either. The data is the foundation; the governance is the safeguard; and the professional judgement of your advisers remains the ultimate quality check.