

PDFテキスト抽出の基本、 問題点、実践的解決策

2020年8月25日 16:00～17:20



アジェンダ

テーマ： PDFテキスト抽出の基本、問題点、実践的解決策

プレゼンター：小林 徳滋 (koba@antenna.co.jp)

時間	見出し	内容
16:00 ～ 16:20	PDFのテキスト抽出の 基本	PDFで文字を表示する仕組み、PDF内のテキストの扱いについて簡単に解説します。
16:20 ～ 16:40	PDFからテキストをコ ピー&ペーストすると きの問題	PDFからテキストをコピー&ペーストする際に経験することがある、様々なトラブルの例を紹介し、トラブルが発生する原因について説明します。
16:40 ～ 17:00	PDFのテキスト抽出の ソリューション	現在、アンテナハウスで開発中のPDF高度テキスト抽出アプリケーションを紹介しま す。簡単なデモを予定しています。
17:00 ～	質疑	チャットでの質疑応答と必要に応じて補足説明します。

※時間は目安です。進行の具合によって多少前後しますのであらかじめご了承ください。



はじめに：作り方によるPDFの大分類

ボーンデジタルPDFと電子化PDF

▶ ボーンデジタルPDF

- ▶ オフィスソフト・DTPソフトなどのアプリで出力したPDF、帳票ソフトやレポート作成ソフトから出力したPDF、など
- ▶ ボーンデジタルPDFは**フォントを使って文字が表示**される

▶ 電子化PDF

- ▶ 紙をスキャンした画像から作るPDF、デジタル写真からPDF変換で作成するPDF
- ▶ 電子化PDFでは**文字はビットマップ画像**で表される。文字の表示にフォントを使わない
- ▶ 電子化PDFのテキスト抽出は、PDFを画像に戻してOCR処理で文字を認識する。OCR文字認識の精度は100%ではないため、誤った文字になる可能性がある

※今日の話題は、主にボーンデジタルPDFからのテキスト抽出に関する内容です。



ボーンデジタルPDFからの テキスト抽出

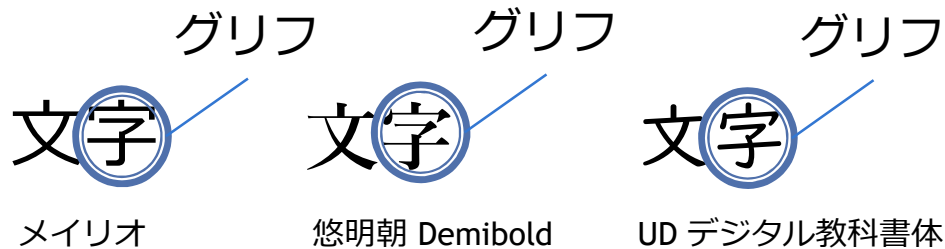
- ▶ たいていのボーンデジタルPDFからは、PDFファイル内にある、文字を表示するためのデータからテキスト抽出できる
こちらが今日のトピックです。
- ▶ そのメリット
 - ▶ OCRと違って文字コードの誤認識がない
 - ▶ 文字の位置を精密に取得できるので位置情報を使える
- ▶ そのデメリット
 - ▶ PDFのセキュリティで、テキストのコピーが許可されている必要がある
 - ▶ 文字コードを取り出せないケースがある
 - ▶ PDFの作り手によって内部が千差万別なため、いろいろなパターンのPDFに対応するための開発工数が大きい



フォント、文字、グリフ

フォントで文字を表示するとは

- ▶ 文字の形（字形）を表す絵のことを**グリフ**という（用語）
- ▶ フォントの中にはグリフを描画するデータがある
- ▶ フォントに文字の識別子を入力して、グリフを探す
- ▶ フォントのエンコーディングは、文字の識別子からフォント内のグリフへの対応関係を示す



※同じ文字でも、フォントによってグリフが異なる

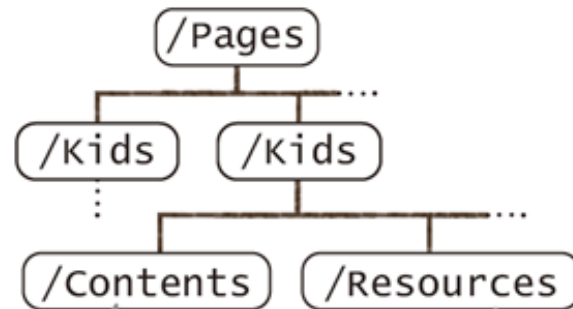


PDFのテキスト表示の仕組み

- ▶ PDFのテキストとはページの特定位置に描かれる文字のこと
- ▶ 表示される文字の識別子と表示位置などのデータは、PDFのページ内のコンテンツストリームに置かれている
- ▶ フォントの情報はフォントリソースに置かれている
- ▶ PDFリーダーは、コンテンツストリーム内の文字の識別子とフォントのエンコーディングを元に、フォントの中のグリフを特定し、指定された位置に文字を印刷・表示する
- ▶ コンテンストリームはページの中で細切れになっている
- ▶ PDF仕様（ISO 32000-1）は文字の識別子からグリフを探して表示する仕組みについて決めている
 - ▶ ISO 32000-1:2008 「第9章 テキスト」を参照



PDFファイル内部の論理構造(一般的な例)



コンテンツストリームが含まれる

```
⋮  
BT  
/F1 11.04 Tf  
1 0 0 1 96.144 716.02 Tm  
[<01260177014A014E>11<0153012A013D>] TJ  
ET  
⋮
```

リソース辞書が含まれる

- フォントや色に関する情報
 - エンコーディング
 - /ToUnicode CMap など

プログラマーから見たPDFファイル

<https://www.antenna.co.jp/pdf/reference/pdftext.html>



PDFの表示

PDF Text Extractor
PDFテキスト抽出




← Times New Roman
← 悠明朝

フォント埋め込みなしで作成したPDFの内部

○コンテンツストリーム

コンテンツストリームに改行コードはない

```
BT
/F0 10.56 Tf 85.08192 729.83997 Td [(P) -1(D) -5(F) -1( ) 34(T) 66(e) 1(xt) 5(
11(E) -2(xt) 5(r) 3(a) 1(c) 1(t) 5(or) ] TJ
/F1 10.56 Tf 0.00004 -17.88 Td [<0050>
4<0044> 4<0046> -250<30C630AD> 12<30B930C8> 11<62BD51FA> ] TJ
ET
```

 使用フォント、フォントサイズ
 行頭の位置
 表示する文字の識別子の並び

○フォントリソース

```
<</Font <</F0 9 0 R/F1 10 0 R>>/
```

```
9 0 obj
```

```
<</Type /Font/Subtype /TrueType/Name /F0/BaseFont
/TimesNewRomanPSMT/FirstChar 32/LastChar 120/Widths [250 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
611 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
500] /Encoding /WinAnsiEncoding/FontDescriptor 11 0 R>>
endobj
```

```
10 0 obj
```

```
<</Type /Font/Subtype /Type0/Name /F1/BaseFont /YuMincho-Regular/Encoding
/UniJIS-UCS2-H/DescendantFonts [12 0 R]>>
endobj
```



PDFの表示

PDF Text Extractor
PDFテキスト抽出

← Times New Roman
← 悠明朝

フォント埋め込みで作成したPDFの内部

コンテンツストリームに改行コードはない

○コンテンツストリーム

```
BT /F0 10.56 Tf 85.08192 729.83997 Td [<0033> -1<0027> -5<0029> -1<0003> 34<0037>  
66<0048> 1<005B0057> 5<0003> 11<0028> -2<005B0057> 5<0055> 3<0044> 1<0046>  
1<0057> 5<00520055> ] TJ  
/F1 10.56 Tf 0.00004 -17.88 Td [<0033> 4<0027> 4<0029> -250<03C403AB>  
12<03B703C6> 11<0BAA095C> ] TJ  
ET
```

■ 使用フォント、フォントサイズ
■ 行頭の位置
■ 表示する文字の識別子の並び

○フォントリソース

```
<</F0 9 0 R/F1 10 0 R>>
```

9 0 obj

```
<</Type /Font/Subtype /Type0/Name /F0/BaseFont  
/JDTVFA+TimesNewRomanPSMT/Encoding /Identity-H/ToUnicode 12 0 R/DescendantFonts  
[11 0 R]>>  
endobj
```

10 0 obj

```
<</Type /Font/Subtype /Type0/Name /F1/BaseFont /JDTVFA+YuMincho-Regular/Encoding  
/Identity-H/ToUnicode 16 0 R/DescendantFonts [15 0 R]>>  
endobj
```



情報処理における文字とテキストデータとは

- ▶ コンピュータのプログラムで文字を取り扱うためには標準化された文字コードを使う。標準文字コードは大雑把には次のようなもの
 - ① 印刷物などで実際に使われている文字を収集する。
 - ② 字義が同じで字体が同一視できる文字を一つにまとめる。これを包摂する（Unification）と言う。文字はこのように抽象化される
 - ③ こうして作成した文字集合の構成要素に一定の番号（コードポイント）を与える。できたものを符号化文字集合と言う
 - ④ コンピュータで情報処理するときはコードポイントを一定の計算式でバイトのならびに変換して取り扱う。これが、情報処理におけるテキストデータであり、この変換方式を符号化方式と言う
- ▶ 符号化文字集合ではUnicodeが一般的。Unicodeをコンピュータで使うときの符号化方式はUTF-8、UTF-16などが使われる
- ▶ 日本の符号化文字集合としてJIS文字規格は幾つかある。JIS X0213:2012が最新のJIS文字規格



PDFテキスト抽出の内部処理

- ▶ PDFでのテキストはコンテンツストリーム内にある文字識別子の並び（表示する文字）のことである
- ▶ 情報処理のテキストは、Unicodeなどで表されたデータ
- ▶ テキスト抽出では、次の二つの処理が必要
 - ① PDF内の文字識別子の並びをUnicodeのテキストデータに変換する（文字コード変換）
 - ② ページ上で文字が表示される位置を考慮しながら、テキストの流れを作る（文の流れの作成）



PDFからテキストを取り出すには

- ▶ ユーザーがPDFファイルからテキストを取り出すとき：
 - ① PDFをリーダーで画面に表示する
 - ② 画面上で、テキストを選択し、コピーする
 - ③ メモ帳などにペーストする
- ▶ ペーストされた結果を保存すると、形式的にはテキストデータになる



PDFテキストの コピー&ペーストの問題例

- ▶ ペーストしたテキストが期待通りになっていないことが多い
- ▶ PDFからテキストをコピー&ペーストしたとき起きる問題の実例のいくつかを紹介



1. 文の順序が表示順と違う

- ▶ 画面に表示している順序とコピー&ペーストしたテキストで段落の順序が異なる
- ▶ Officeソフトから作成したシンプルなレイアウトのPDFでも起きることがある（例1）
- ▶ DTPソフトで作成した複雑なレイアウトのPDFではかなり頻繁にみられる（例2）



例 1

休業実績一覧表（様式特小第2号（小規模事業主用様式））

【記入要領】

1 「支給申請する1か月間（判定基礎期間）」

雇用調整助成金では、原則として1か月単位で休業の実績について確認し、それに基づいて支給がなされます。この休業の実績を判定する1か月単位の期間を「判定基礎期間」といいます。この単位になる1か月は、休業した事業所の毎月の賞金の締め切り日の翌日から、その次の締め切り日までの期間です。毎月の賞金の締め切り日がない場合などは、雇月（カレンダーの1日～月末日）となります。

支給申請のときは、この判定基礎期間を単位として、令和2年1月24日以降の期間であれば複数月分をまとめて申請することができます。その場合、この「休業実績一覧表」は、判定基礎期間ごとに作成する必要があります。

2 「従業員の数」

2か月を超えて使用される者（実働として2か月を超えて使用されている者のほか、それ以外の者であっても雇用期間の定めのない者及び2か月を超える雇用期間の定めのある者を含む。）であり、かつ、適当な所定労働時間が、当該事業主に雇用される通常の労働者と概ね同等（現に当該事業主に雇用される通常の労働者の適当な所定労働時間が40時間である場合は、概ね40時間である者をいう。ただし、労働基準法（昭和22年法律第49号）の特例として、所定労働時間がまだ40時間を上回っている場合は、「概ね同等」とは、概ね当該所定労働時間を指す。）である者をいいます。

3 「休業手当支払い率」

労働者の代替と休業の方法について約率するときに決めた休業手当の支払い率を、記入してください。

なお、休業するときに労働者に支払う休業手当の額は、通常支払っている賞金の60%以上である必要があります。通常支払っている賞金と同じ額を支払っているときは、100%と記入してください。

対象労働者ごとにちがう複数の支払い率がある場合は、最も多い労働者に適用している支払い率としてください。または、すべての支払い率の単純平均か、加重平均で計算した支払い率でもかまいません。

（例）支払い率が60%の従業員5人、80%の従業員2人、100%の従業員3人の場合

最も多い従業員に適用している支払い率：60%

単純平均： $(60 \times 5 + 80 \times 2 + 100 \times 3) \div 3$ 種類=80（%）

加重平均： $(60 \times 5 + 80 \times 2 + 100 \times 3) \div 10$ 人=76（%）のうちいずれかを選択

丸一日休業した場合と、1日のうちの一部休業した場合で、ちがう支払い率としている場合は、加重平均で計算した支払い率としてください。

（例）丸一日休業したときの支払い率90%で10日休業し、

1日のうち一部休業したときの支払い率80%で3日分（※）休業した場合 ※4で一日換算した⑤の日数

加重平均： $(90 \times 10 + 80 \times 3) \div 13$ 日=88（%）

4 休業した事業所で、労働者が通常1日に働く労働時間数を記入してください。（就業規則や雇用契約書、労働条件通知書などに記載している労働時間数です。）

労働者ごとに所定労働時間がちがう場合は、最も多い労働者に適用している所定労働時間数としてください。

5 ⑤欄は「④一部の時間帯休業した時間数」の合計を、4の時間数で割ることにより、何日分休業したことになるかを計算して、記入してください。小数点以下は切り上げて、整数で記入してください。

6 休業対象労働者ごとの休業実績一覧

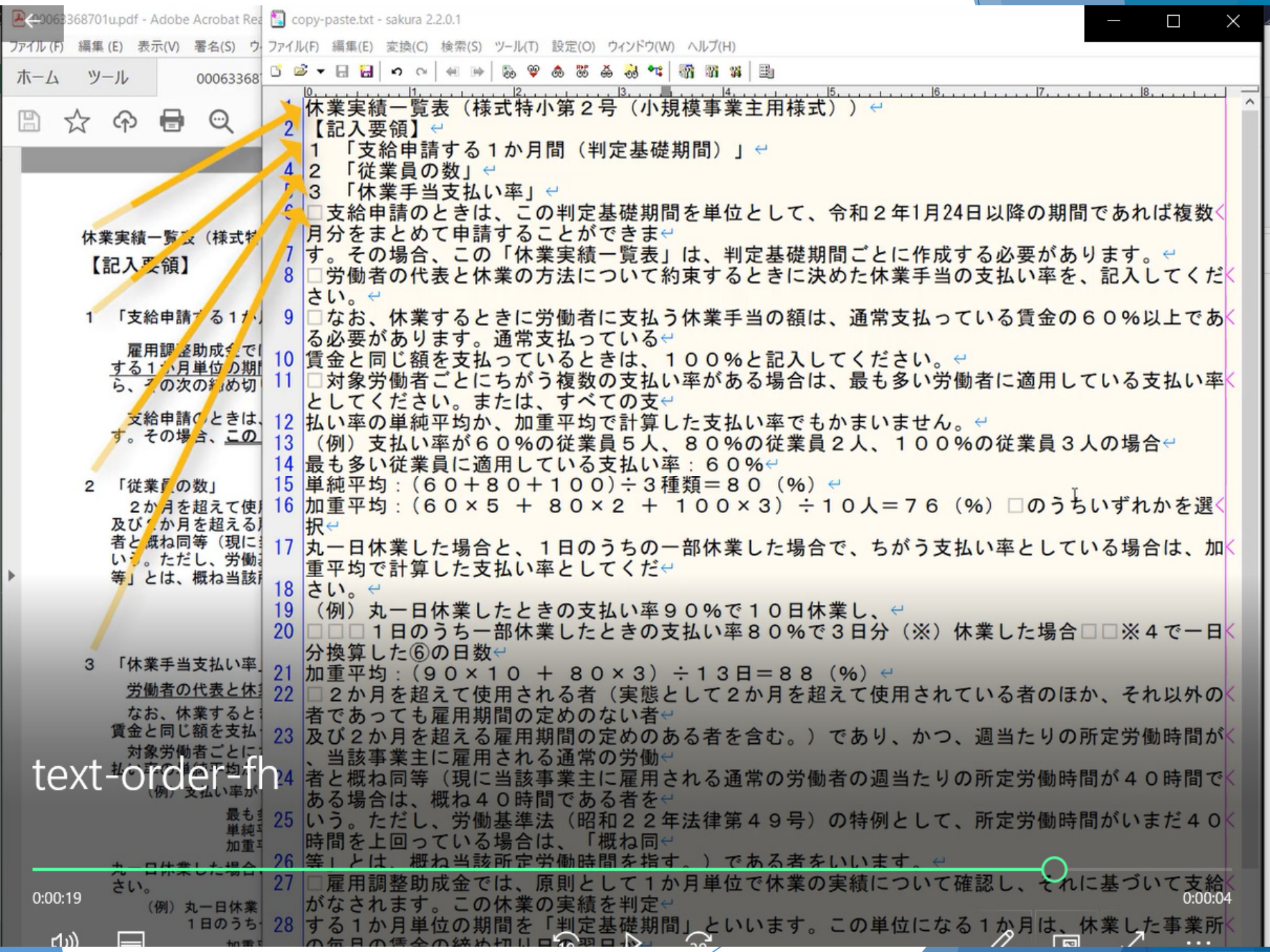
解雇予告をされた者、退職願を提出した者、退職勧奨に応じた者は含めることができません（それぞれの時点より前までの休業についてであれば含めて可）。また、供給ができない他の助成金の対象労働者は含めることができません。

- ① 氏名 できれば、添付して提出する「休業手当の額がわかる書類（賞金台帳や給与明細など）」か「休業させた日や時間がわかる書類（出勤簿やタイムカードなど）」と同じ順番になるように記入してください。
- ② 雇用保険被保険者番号 対象労働者を雇用保険に加入させたときに発行される「雇用保険被保険者証」などに記載されている被保険者番号を記入してください。（※雇用保険に加入させていない従業員については、「緊急雇用安定助成金」というもうひとつの助成金の申請をしてください。）
- ③ 1日休業した日数 丸一日休業した日数を記入してください。
- ④ 1日のうち一部休業した時間数 通常の一日の営業時間のうち、一部休業した場合の時間数を記入してください。具体的には、短時間休業の時間（30分未満は切り捨て、例：1時間40分→1.5）数の合計を記入してください（合計欄は小数点以下切り上げ。）。
- ⑤ 判定基礎期間の休業手当の額 「判定基礎期間」中の休業について、「休業手当支払い率」にもとづき、対象労働者ごとに支払う休業手当の額を記入してください。1日休業した場合（③）と1日のうち一部休業した時間数（④）の額をわける必要はありません。
※③～⑤欄は、一覧に記載した対象労働者すべての数を合計して、合計欄に記入してください。
- ⑦ 休業延べ日数 ③の合計（一覧に記載した対象労働者ごとの丸一日休業した日数の合計）と、⑤の日数の、合計日数を記入してください。

7 一覧表の下にある文章を確認し、記名押印または署名してください。また、休業に関する内容が事前に確約した内容であることについて労働者代表の方に確認してもらい、記名押印または署名をもらってください。

様式特小第2号
（新型コロナウイルス感
染症関係）
（小規模事業主用様式）
（厚生労働省）





休業実績一覧表 (様式特小第2号 (小規模事業主用様式))

【記入要領】

- 1 「支給申請する1か月間 (判定基礎期間)」
- 2 「従業員の数」
- 3 「休業手当支払い率」

支給申請のときは、この判定基礎期間を単位として、令和2年1月24日以降の期間であれば複数

月分をまとめて申請することができます。その場合、この「休業実績一覧表」は、判定基礎期間ごとに作成する必要があります。

労働者の代表と休業の方法について約束するときには決めた休業手当の支払い率を、記入してください。

なお、休業するとき労働者に支払う休業手当の額は、通常支払っている賃金の60%以上である必要があります。通常支払っている

賃金と同じ額を支払っているときは、100%と記入してください。

対象労働者ごとにちがう複数の支払い率がある場合は、最も多い労働者に適用している支払い率としてください。または、すべての支

払い率の単純平均か、加重平均で計算した支払い率でもかまいません。

(例) 支払い率が60%の従業員5人、80%の従業員2人、100%の従業員3人の場合

最も多い従業員に適用している支払い率: 60%

単純平均: $(60 + 80 + 100) \div 3 \text{種類} = 80 (\%)$

加重平均: $(60 \times 5 + 80 \times 2 + 100 \times 3) \div 10 \text{人} = 76 (\%)$ □のうちいずれかを選択

丸一日休業した場合と、1日のうちの一部休業した場合で、ちがう支払い率としている場合は、加

重平均で計算した支払い率としてください。

(例) 丸一日休業したときの支払い率90%で10日休業し、

□□□1日のうち一部休業したときの支払い率80%で3日分 (※) 休業した場合 □□※4で一日

分換算した⑥の日数

加重平均: $(90 \times 10 + 80 \times 3) \div 13 \text{日} = 88 (\%)$

2か月を超えて使用される者 (実態として2か月を超えて使用されている者のほか、それ以外の

者であっても雇用期間の定めのない者

及び2か月を超える雇用期間の定めのある者を含む。) であり、かつ、週当たりの所定労働時間が

、当該事業主に雇用される通常の労働者

者と概ね同等 (現に当該事業主に雇用される通常の労働者の週当たりの所定労働時間が40時間

ある場合は、概ね40時間である者を

いう。ただし、労働基準法 (昭和22年法律第49号) の特例として、所定労働時間がいまだ40

時間を上回っている場合は、「概ね同

等」とは、概ね当該所定労働時間を指す。) である者をいいます。

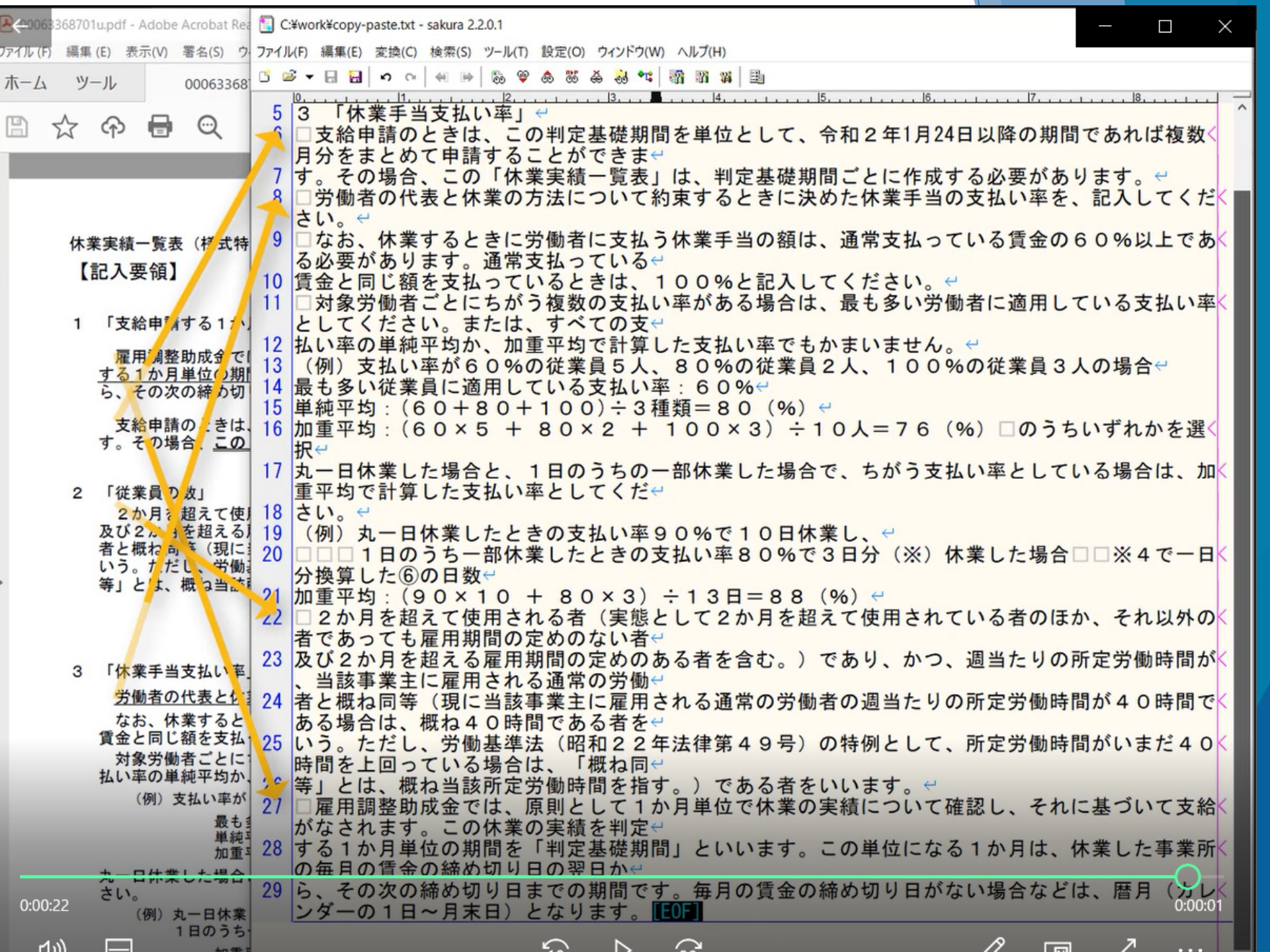
雇用調整助成金では、原則として1か月単位で休業の実績について確認し、それに基づいて支給

がなされます。この休業の実績を判定

する1か月単位の期間を「判定基礎期間」といいます。この単位になる1か月は、休業した事業所

の毎月の賃金の締め切り日(翌日)

text-order-fh



3 「休業手当支払い率」

□支給申請のときは、この判定基礎期間を単位として、令和2年1月24日以降の期間であれば複数ヶ月分をまとめて申請することができます。

す。その場合、この「休業実績一覧表」は、判定基礎期間ごとに作成する必要があります。
□労働者の代表と休業の方法について約束するときに決めた休業手当の支払い率を、記入してください。

□なお、休業するときに労働者に支払う休業手当の額は、通常支払っている賃金の60%以上である必要があります。通常支払っている賃金と同じ額を支払っているときは、100%と記入してください。

□対象労働者ごとにちがう複数の支払い率がある場合は、最も多い労働者に適用している支払い率としてください。または、すべての支払い率の単純平均か、加重平均で計算した支払い率でもかまいません。

(例) 支払い率が60%の従業員5人、80%の従業員2人、100%の従業員3人の場合
最も多い従業員に適用している支払い率：60%

単純平均：(60+80+100)÷3種類=80(%)
加重平均：(60×5+80×2+100×3)÷10人=76(%) □のうちいずれかを選択

丸一日休業した場合と、1日のうちの一部休業した場合で、ちがう支払い率としている場合は、加重平均で計算した支払い率としてください。

(例) 丸一日休業したときの支払い率90%で10日休業し、
□□□1日のうち一部休業したときの支払い率80%で3日分(※)休業した場合□□※4で一日分換算した⑥の日数

加重平均：(90×10+80×3)÷13日=88(%)
□2か月を超えて使用される者(実態として2か月を超えて使用されている者のほか、それ以外の者であっても雇用期間の定めのない者

及び2か月を超える雇用期間の定めのある者を含む。)であり、かつ、週当たりの所定労働時間が、当該事業主に雇用される通常の労働者と概ね同等(現に当該事業主に雇用される通常の労働者の週当たりの所定労働時間が40時間である場合は、概ね40時間である者を

いう。ただし、労働基準法(昭和22年法律第49号)の特例として、所定労働時間がいまだ40時間を上回っている場合は、「概ね同等」とは、概ね当該所定労働時間を指す。)である者をいいます。

□雇用調整助成金では、原則として1か月単位で休業の実績について確認し、それに基づいて支給がなされます。この休業の実績を判定する1か月単位の期間を「判定基礎期間」といいます。この単位になる1か月は、休業した事業所の毎月の賃金の締め切り日の翌日から、その次の締め切り日までの期間です。毎月の賃金の締め切り日がない場合などは、暦月(カレンダーの1日~月末日)となります。 [EOF]

ら、その次の締め切り日までの期間です。毎月の賃金の締め切り日がない場合などは、暦月(カレンダーの1日~月末日)となります。 [EOF]

ら、その次の締め切り日までの期間です。毎月の賃金の締め切り日がない場合などは、暦月(カレンダーの1日~月末日)となります。 [EOF]

休業実績一覧表 (様式特) 【記入要領】

1 「支給申請する1か月の期間」

雇用調整助成金で支給申請する1か月単位の期間から、その次の締め切り日の前日までの期間を1か月の期間としてください。

支給申請のときは、この判定基礎期間を単位として、令和2年1月24日以降の期間であれば複数ヶ月分をまとめて申請することができます。

2 「従業員の数」

2か月を超えて使用される者(実態として2か月を超えて使用されている者のほか、それ以外の者であっても雇用期間の定めのない者)及び2か月を超える雇用期間の定めのある者を含む。)であり、かつ、週当たりの所定労働時間が、当該事業主に雇用される通常の労働者と概ね同等(現に当該事業主に雇用される通常の労働者の週当たりの所定労働時間が40時間である場合は、概ね40時間である者をいう。ただし、労働基準法(昭和22年法律第49号)の特例として、所定労働時間がいまだ40時間を上回っている場合は、「概ね同等」とは、概ね当該所定労働時間を指す。)である者をいいます。

3 「休業手当支払い率」

労働者の代表と休業の方法について約束するときに決めた休業手当の支払い率を、記入してください。
なお、休業するときに労働者に支払う休業手当の額は、通常支払っている賃金の60%以上である必要があります。通常支払っている賃金と同じ額を支払っているときは、100%と記入してください。
対象労働者ごとにちがう複数の支払い率がある場合は、最も多い労働者に適用している支払い率としてください。または、すべての支払い率の単純平均か、加重平均で計算した支払い率でもかまいません。

(例) 支払い率が60%の従業員5人、80%の従業員2人、100%の従業員3人の場合
最も多い従業員に適用している支払い率：60%

単純平均：(60+80+100)÷3種類=80(%)
加重平均：(60×5+80×2+100×3)÷10人=76(%) □のうちいずれかを選択

特集 弾ける美味さ！
信州中野のおいしいぶどう



武田政志 さん（竹原）



日本一の産地を目指して！ 味にこだわる信州中野のぶどう

わが家のぶどう栽培

両親と妻と私の4人で「巨峰」「ナガノパール」「ピオーネ」「シャインマスカット」「クインシーナ」など8種類のぶどうを栽培しています。わが家は元々リンゴなどを栽培する農家でしたが、祖父の代からぶどうの栽培を始め、私はぶどう農家として3代目になります。

顔を見ると嬉しいもので、来年も良いものを作らなければという意欲が湧いてきます。

中野市のぶどうの歴史
本市で本格的に巨峰の栽培が始まったのは昭和28年ごろだといわれています。ぶどうといえば「アラウエア」や「ナイアガラ」が主流だった当時、大粒で味が抜群に良い「巨峰」が登場しました。当初は栽培が難しく安定した生産ができなかったようですが、次第に栽培技術が発達し、安定生産が可能になったことで栽培農家が増え、中野市は日本有数の

の巨峰の産地となりました。現在では、5月に出荷する超加温栽培、6～8月に出荷する加温・無加温栽培、9～11月に出荷する露地栽培、12月～1月に出荷する冷蔵出荷と、5月から翌年1月までの長期間にわたり、味にこだわったぶどうを出荷しています。

仲間との支え合い

地域のぶどう農家には、産地の歴史を築いてきた先輩方だけでなく、若い後継者が入ってきて活気があり、とても心強く感じています。

畑で栽培技術の話をしたり、お互いのハウスのビニールを張る作業をしたりして、仲間同士支え合ってぶどう栽培をしています。昨年のことですが、突風でハウスのビニールが捲れ上がってしまったときに、園主が不在にもかかわらず近所の仲間たちが集まって直してくれたという出来事もありました。

中野のぶどうを全国へ
地元で採れた新鮮でおいしいぶどうを「おつかいもの」などに使っていたとき、「信州中野のぶどう」を自慢の逸品として皆さんにも宣伝していただければと思います。

信州中野発！

おいしいぶどうがとれるまで

甘くて果汁がたっぷりのおいしいぶどうはどのよう
に栽培されている
のでしょうか？
巨峰の露地栽培
の様子を、JA中
野市の営農指導員
の方に聞きました。

おいしいぶどうは、
農家の皆さんの努力と
愛情の結晶です！



▲JA中野市でぶどうの営農指導員を務める清水達哉さん



②発芽（4月）：春を迎え、ぶどうの木から若い新芽が顔を出します



①剪定（12月～2月上旬）：枝を切ることでより降雪によるぶどう棚の倒壊を防ぎます



⑨着色（8月）：巨峰は1粒ずつ着色が進んでいきます



⑩収穫（9月）：手間暇かけて育てたぶどうがついに収穫を迎えます

中野市
別冊 広報なかの vol.1

P. 10
https://www.city.nakano.nagano.jp/docs/2015121800024/file_contents/1011.pdf





特集
弾ける美味さ！
信州中野のおいしいぶどう

わが家のぶどう栽培

両親と妻と私の4人で「巨峰」、「ナガノパープル」、「ピオーネ」、「シャインマスカット」、「クイーンニーナ」など8種類のぶどうを栽培しています。わが家は元ナガノパープルなどを栽培する農家でしたが、祖父の代からぶどうの栽培を始め、私はぶどう農家としては3代目になります。

現在は「JA中野市ぶどう部会」の部会長を務めており、県外へぶどうのセールスに行く機会も多くあります。お客様の「おいしい」という笑顔を見ると嬉しいもので、来年も良いものを作らなければという意欲が湧いてきます。

中野市のぶどうの歴史

本市で本格的に巨峰の栽培が始まったのは昭和28年ごろだといわれています。ぶどうといえば「デラウェア」や「ナイアガラ」が主流だった当時、大粒で味が抜群に良い「巨峰」が登場しました。当初は栽培が難しく安定した生産ができなかったようですが、次第に栽培技術が発達し、安定生産が可能になったことで栽培農家が増え、中野市は日本有数の巨峰の産地となりました。

現在では、5月に出荷する超加温栽培、6～8月に出荷する加温・無加温栽培、9～11月に出荷する露地栽培、11月に出荷する抑制栽培、12～1月に出荷する冷蔵出荷と、5月から翌年1月までの長期間にわたり、味にこだわったぶどうを出荷しています。

仲間との支え合い

地域のぶどう農家には、産地の歴史を築いてきた先輩方だけでなく、若い後継者が入ってきて活気があり、とても心強く感じています。

畑で栽培技術の話をしたり、お互いのハウスのビニールを張る作業をしたりして、仲間同士支え合ってぶどう栽培をしています。昨年のごとく、突風でハウスのビニールが捲れ上がってしまったときに、園主が不在にもかかわらず近所の仲間たちが集まって直してくれたという出来事もありました。

中野のぶどうを全国へ

地元で採れた新鮮でおいしいぶどうを「おつかいもの」などに使っていただき、「信州中野のぶどう」を自慢の逸品として皆さんにも宣伝していただければと思います。

日本一の産地を目指してー。
味にこだわる信州中野のぶどう

武ただだ田政まさし志
 (竹原)
 さん

信州中野発！

甘くて果汁がたっぷりのおいしいぶどうはどのように栽培されているのでしょうか？
 巨峰の露地栽培の様子を、JA中野市の営農指導員の方に聞きました。

おいしいぶどうができるまで

- ②発芽（4月）：春を迎え、ぶどうの木から若い新芽が顔を出します
- ①剪定（12月～2月上旬）：枝を切るにより降雪によるぶどう棚の倒壊を防ぎます
- おいしいぶどうは、農家の皆さんの努力と愛情の結晶です！**
- ⑩収穫（9月）：手間暇かけて育てたぶどうがついに収穫を迎えます
- ⑨着色（8月）：巨峰は1粒ずつ着色が進んでいきます
- ▲JA中野市でぶどうの営農指導員を務める清しみず水達たつや哉さん

10
別冊
なかの
広報
vol.1



2. 文の配置とテキストの順序がうまく対応しない

- ▶ コピー&ペーストで、文字や文の配置とテキストの順序が不適切になる
 - ▶ 次のスライドに定義箇条形式からのテキスト抽出の例をしめす。
- ▶ 2段組の2つの段でテキストがつながってしまうことがある
- ▶ 表や図のコピー&ペーストも面倒
 - ▶ 表や図はむしろ画像にしてしまう方が良いときもある



Title	<i>DITA Open Toolkit User Guide and Reference.</i>
Edition, release	First edition, August 10, 2006. Based on release 1.2.2 of DITA Open Toolkit.
Publishing information	DITA Open Toolkit is an open source, reference implementation of the OASIS DITA standard (currently DITA 1.0)
Authors	Anna van Raaphorst and Richard H. (Dick) Johnson, principals, VR Communications, Inc. (http://www.vrcommunications.com).
Description	This document is the definitive source of information about DITA Open Toolkit (OT). It is also a product of the architecture and the recommended best practices, having been written entirely in DITA XML and produced using the principles and procedures described in the document. With a few minor exceptions (in anticipation of some major enhancements in release 1.3, we made use of the existing bookmap and PDF2 functionality), <i>DITA Open Toolkit User Guide and Reference</i> is "vanilla DITA." By processing this document to a given target environment you can see output with no specializations or other special configurations applied. In general, we have found the vanilla outputs adequate for our needs, if not always glamorous or exciting.

※定義と説明文が入り組む ↓ テキストエディタに貼り付け

Title DITA Open Toolkit User Guide and Reference.

Edition, release First edition, August 10, 2006. Based on release 1.2.2 of DITA Open Toolkit.

DITA Open Toolkit is an open source, reference implementation of the OASIS DITA standard (currently DITA 1.0)

Publishing information

Anna van Raaphorst and Richard H. (Dick) Johnson, principals, VR Communications, Inc. (<http://www.vrcommunications.com>).

Authors

This document is the definitive source of information about DITA Open Toolkit (OT). It is also a product of the architecture and the recommended best practices, having been

Description

written entirely in DITA XML and produced using the principles and procedures described in the document. With a few minor exceptions (in anticipation of some major enhancements in release 1.3, we made use of the existing bookmap and PDF2 functionality), *DITA Open Toolkit User Guide and Reference* is "vanilla DITA." By processing this document to a given target environment you can see output with no specializations or other special configurations applied. In general, we have found the vanilla outputs adequate for our needs, if not always glamorous or exciting.



3. 文字間の空き、改行の扱い

- ▶ 和欧混色の場合、元のテキストにはスペースがないが、コピー&ペースト後にスペースが入る
- ▶ 行末の改行の扱い

※コピー&ペーストを処理するPDFリーダーによる相違がある



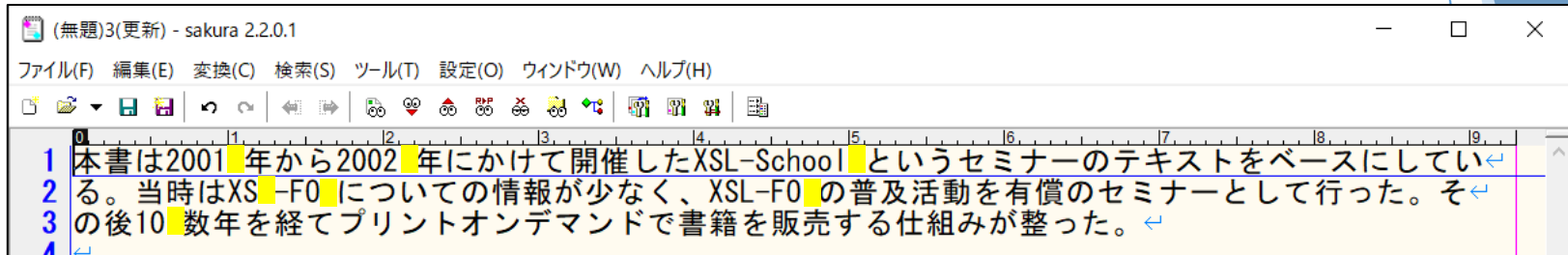
Adobe Reader で選択して、コピー



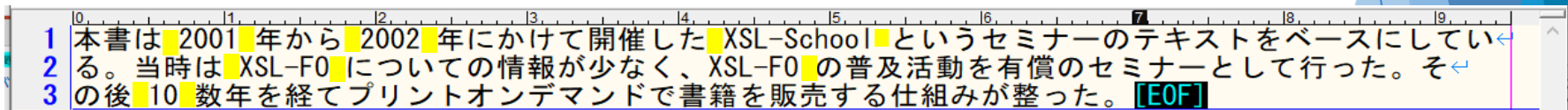
半角空白が挿入される（黄色のところ）



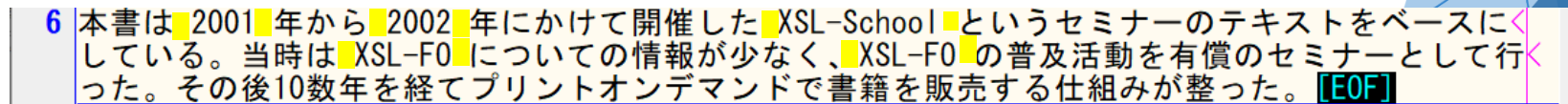
テキストエディタに貼り付け



Google Chromeでコピーすると、半角空白挿入位置が異なる



Firefoxでコピーすると、さらに結果が異なる。行末に改行が入らない。



PDF内のテキストには和文と欧文間の空白文字はない。PDFの表示オペレータ (Tj/TJ) のパラメータで表示位置の調整が行われている。改行も同様。



4. 不要な文字が抽出される

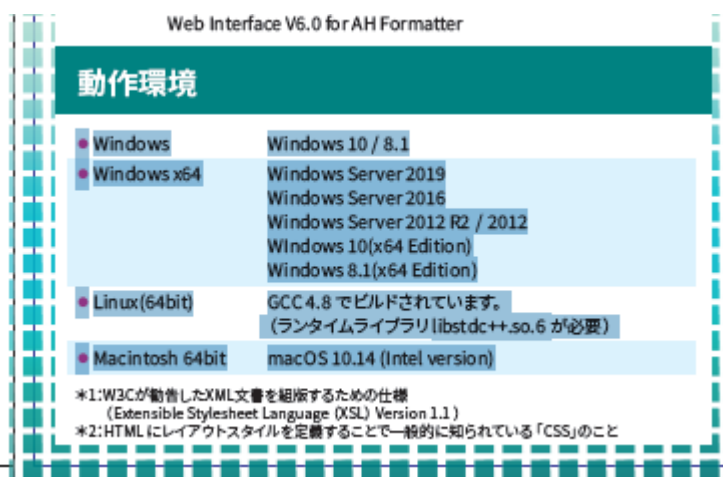
- ▶ 文字が2重に出る ⇒ 強調を文字の2重表示で表現しているとき、コピー&ペーストすると文字が2重になる。
最近のAdobe Readerは余り起きない
- ▶ 表示されていない文字が取れる
⇒ PDFで非表示に設定されている文字など
- ▶ 不必要なテキストがコピーされる
 - ▶ ヘッダー・フッター（柱）は抽出不要なので無視したい
 - ▶ ページ番号（ノンブル）は抽出不要なので無視したい



5. 文字コードを取得できない

※この問題には回避策がないことが多い。

- ▶ PDFのテキストを正しい標準文字コードに変換できない
 - ▶ カスタムエンコーディング、ユーザー定義文字などが使われている
- ▶ グリフがアウトライン化（画像化）されている



Adobe Acrobat Readerで選択、コピー
テキストエディタにペースト⇒

※ラテン文字が抽出できない

でビルドされています。
(ランタイムライブラリ が必要)



コピー＆ペーストには問題多い

- ▶ **コピー＆ペーストはPDF内のテキストを再利用するための十分な解決策にならない！**
- ▶ その原因は、PDFファイルの作成目的は、テキストを**表示・印刷**することにあるため
- ▶ **新しいツールが必要！**
- ▶ **簡単そうでも困難な課題**



開発中

テキスト抽出ツール

PDF Advanced Extractor α版

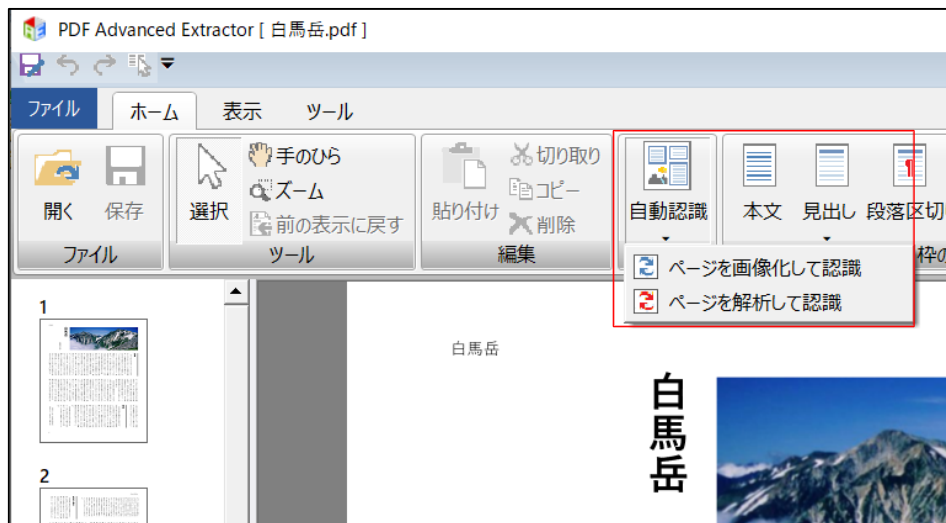
- ▶ PDFファイルからテキストデータを抽出する
- ▶ PDFのテキストをコピー&ペーストで起きる典型的な問題として示した、例1～4については本ツールを使えば解決できる

※ PDFファイルを解析する方法のため、残念ながら、例5は解決できない



テキスト枠を自動認識

▶ 2段組、段抜きなどの領域を自動認識



1 白馬岳

2 白馬岳

3 今では日本北アルプスの名で広く世に知られている飛騨山脈は、加藤厚平士の説に拠ると、丹波山脈より北十度東に向って並走する数条の連脈から成っているものであるという、其連脈の一に白馬山脈というのがある。立山山脈との対称上から又後立山山脈とも呼ばれ、飛騨山脈中の最も長い山脈で、北は日本海岸の奥不知附近から起り、越中と越後及び信濃との国境を測走して遠く飛騨国内に達しているが、中に於て越中、越後及び信濃の三国界から飛騨、信濃及び越中の三国界附近に至る、直径にして五十五并約四里の間が主要部ともいふ可き部分であつて、最高二千九百九十米、最低二千八百八十米、平均高度は二千六百二十米に及んでゐる。そして二千八百米を超えてゐる峰は十五、六座を下らないのである。それが松本市の西の縁から大原を建てたように急に聳え立っているので、地形の相違の著しい為、二千五百米以下には中山性の地勢と称す可きものに属するに拘わらず、恰も大山脈を見るが如き観を呈し、加うるに盛夏八月の候も尚白皚白に輝く雪田が山の嶺を飾り、雪深が幾まとなぐ山肌に象眼されているので、頂上附近の高山性地勢と稱して、一層崇峻雄大な感じを起さしめるのである。

4 白馬山脈の最高峰は、中央より稍や南に偏している東岳であつて、水晶や新水晶などを産する所から水晶山の名もある。三尖点の位置は駒形より十米余も低い峰に在るので、其の高さは恐らく二千九百九十米を下ることはあるまい。之に次ぐものは主要部の北端に在る白馬岳で、海拔高度二千九百三十三米、最高点は長野県北安曇郡と富山県下新川郡に跨り、支脈北に向つて行くこと十町余り山脈は二岐し、其間に新潟県平野郡を包んでいる。越中、越後に言つると、白馬岳は一部分しか新潟県には跨っていないことになる。松本市から越後の糸魚川町に達する糸魚川街道は、平地から此山脈を仰望するに最も適した街道であつて、五月下旬、麓の新雪が薄く濃やかならんとする

5 白馬山脈を構成する岩石は、大体に於て花崗岩又は之に類似した深成岩であるが、新大岩が其間に噴出し、又古生層の露出せる所も少なくない。然し調査が充分に行き届いていないので、はなから、詳細に探究された際には、新に発見する所が少なくないであらうと思ふ。白馬岳の近傍は花崗岩、蛇紋岩、古生層、礫岩及び新火山岩などで構成されている。そして頂上附近は概して東側は古生層、西側は礫岩から成り、中央の一部は蛇紋岩の認められることが、地質調査所の地図に明示されている。

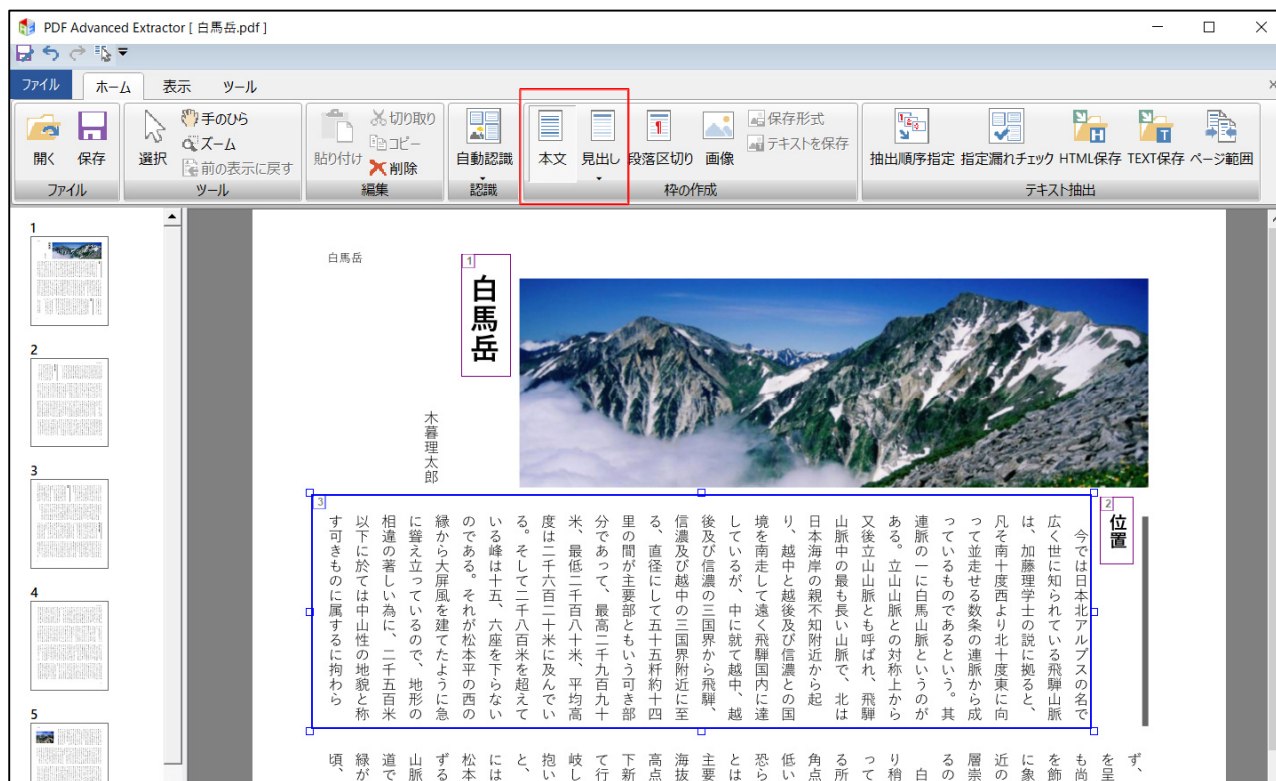
6 白馬岳の南には杓子岳があり、更に其南に據つて鑓ヶ岳がある。假に之を白馬三山と稱え、共に同じ地質から成つてい

7 昭和十二年二月東京地誌第10巻刊の「白馬岳」図説に拠

8 い雪山の姿を望むことは、我々の山岳遊覧中に在りて優れたるの山の一つであるといつてはよい。

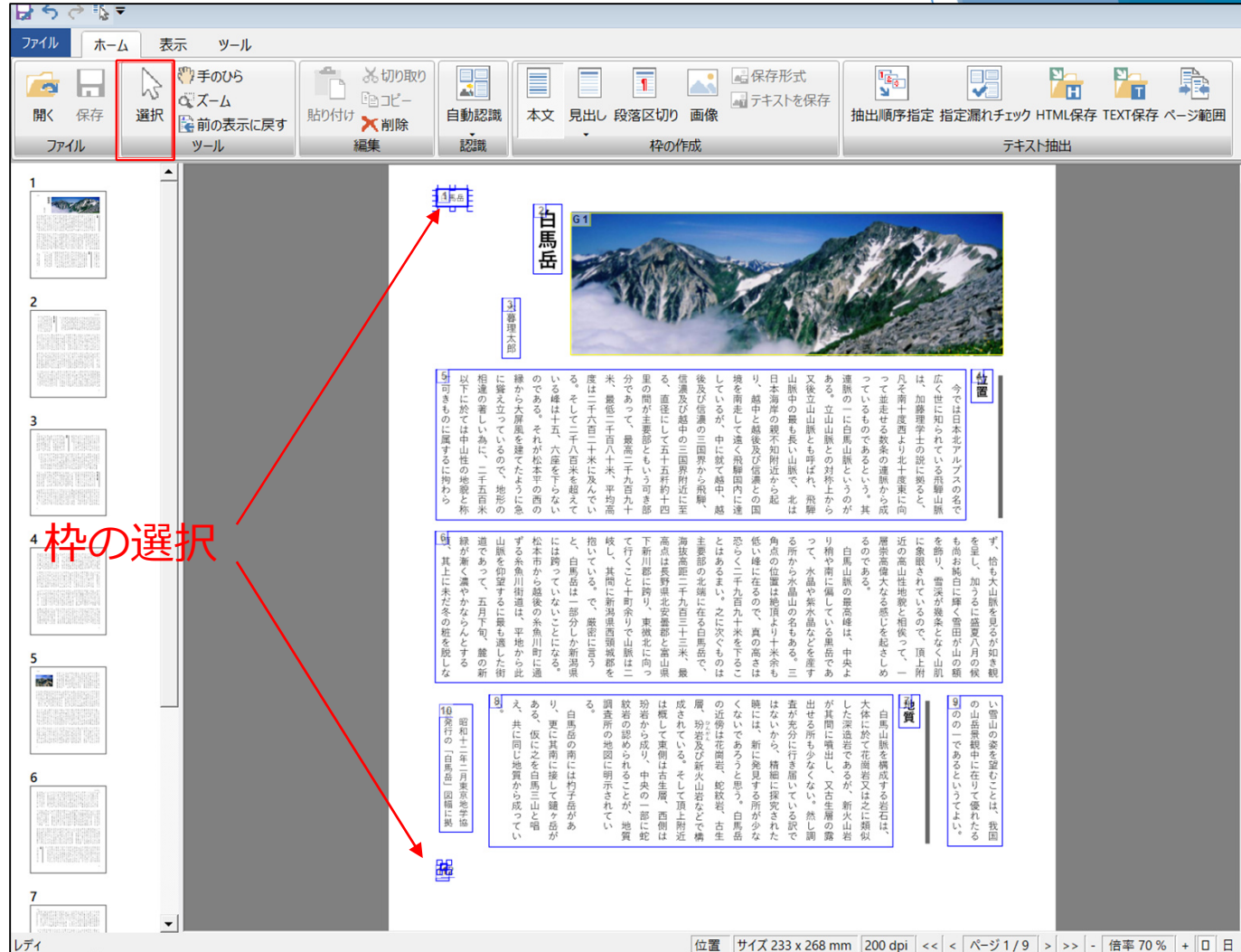
テキスト枠を作成

- ▶ 画面上でマウスをドラッグしてテキスト枠を指定
- ▶ 見出しの指定もできる



テキスト枠の編集

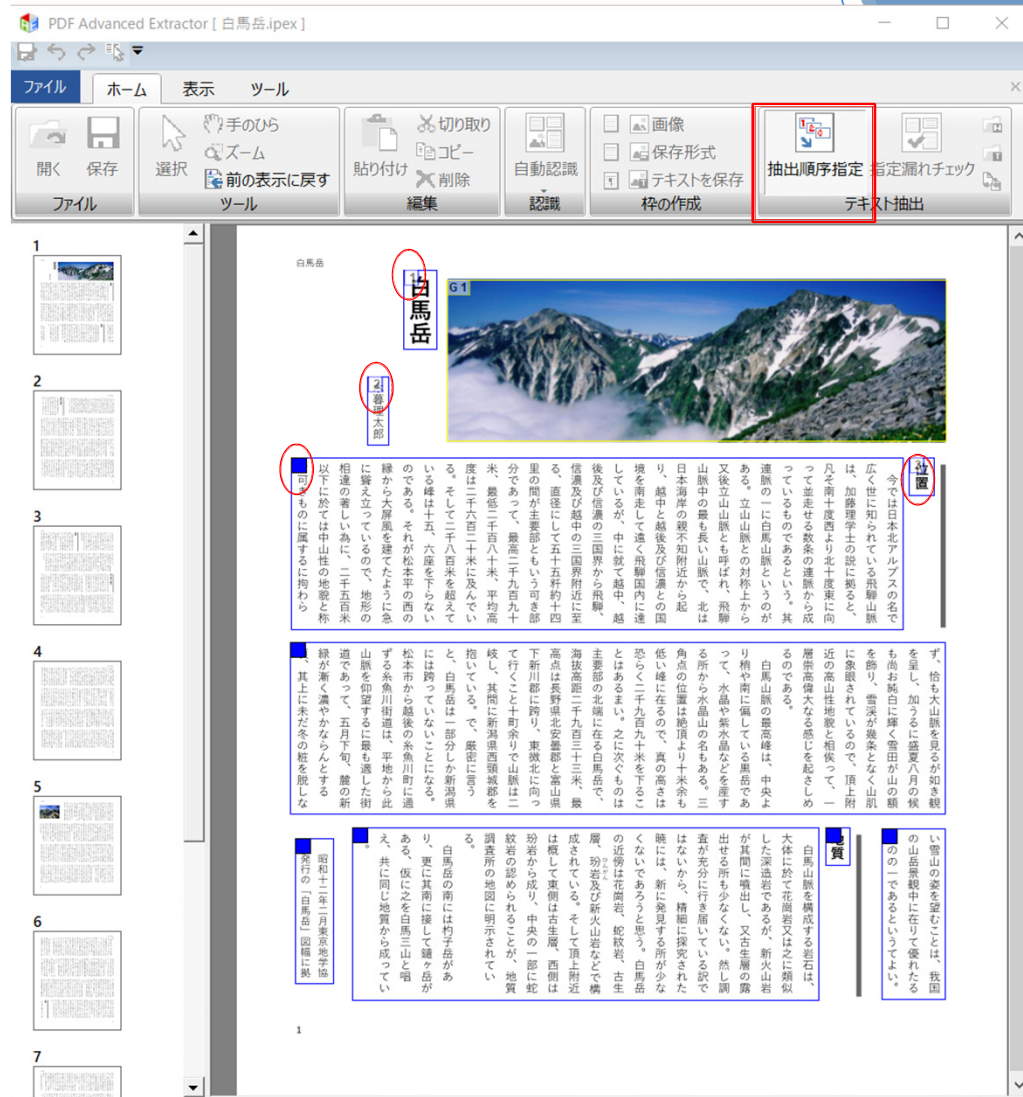
- ▶ 枠選択
- ▶ 枠拡大
- ▶ 枠縮小
- ▶ 枠移動
- ▶ 枠結合
- ▶ 枠削除



枠の選択

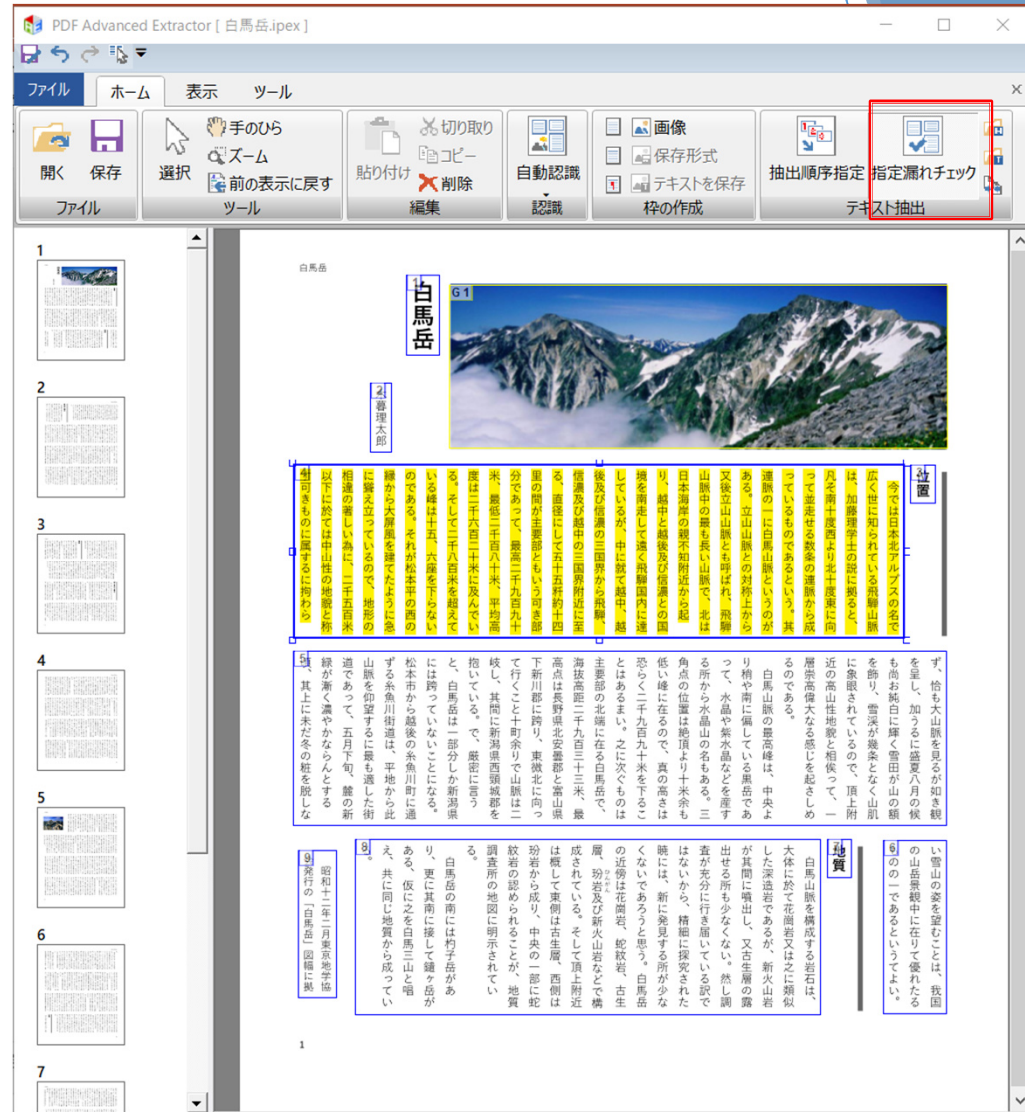
テキスト枠の抽出順序変更

- ▶ テキスト枠の内容を出力する順序を指定
- ▶ 既定値は作成順。必要に応じて付け替え



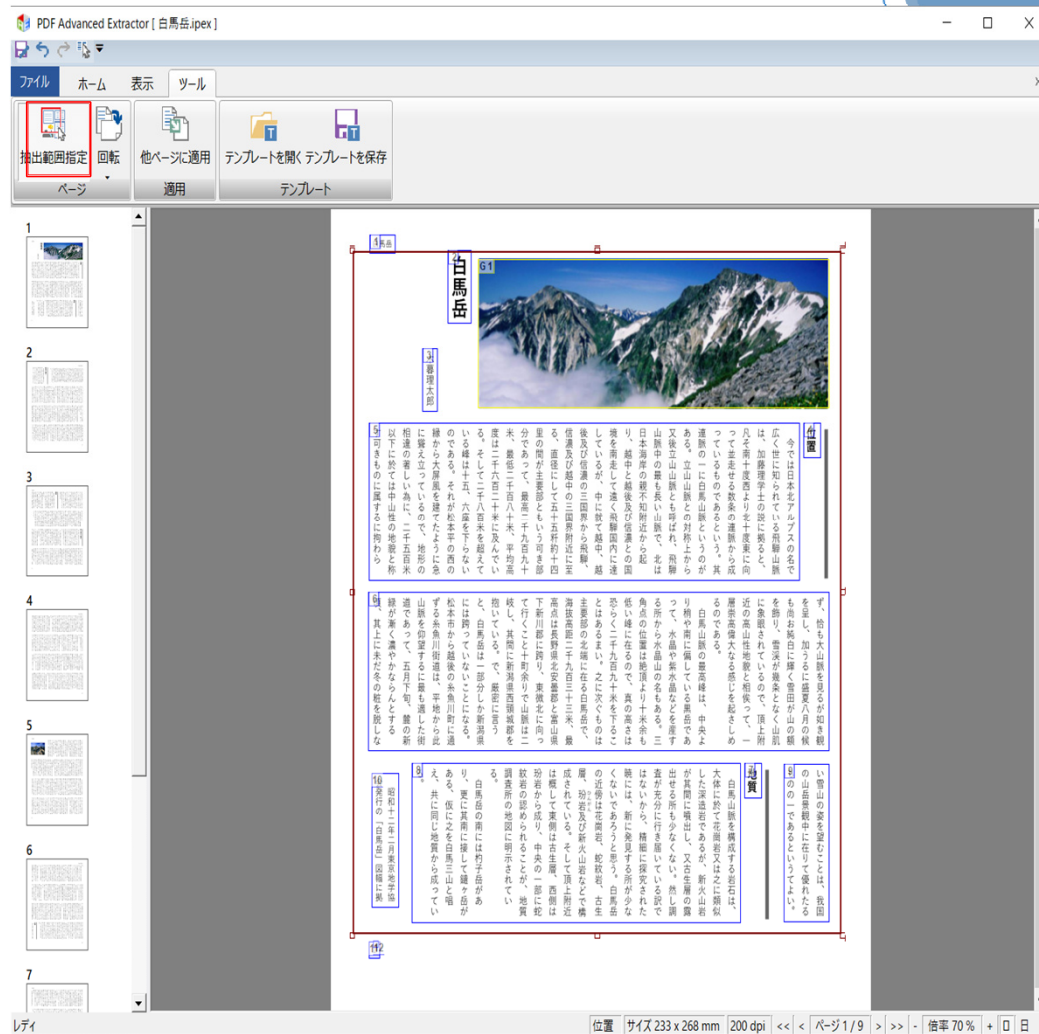
指定漏れチェック

- ▶ テキスト枠毎に、文字を抽出したときに脱落しないかを事前に確認できる



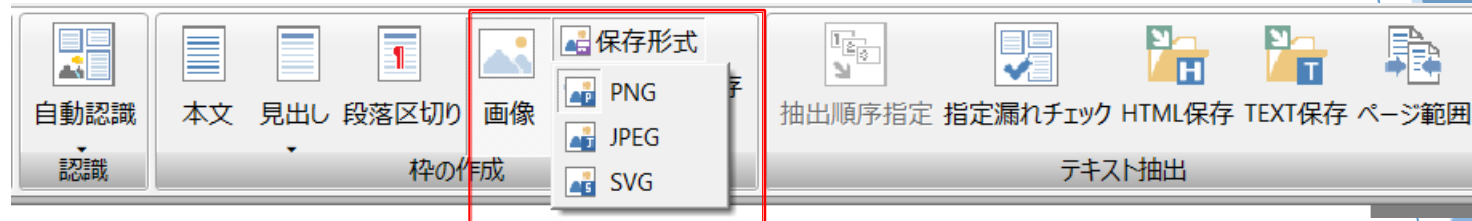
抽出範囲指定

- ▶ ページの中でテキストを抽出する範囲を指定する
- ▶ 本文領域を範囲指定すると、ヘッダー、フッターを除外できる



画像保存

- ▶ 画像領域を指定して画像ファイルを作成



- ▶ 選択範囲を画像化。表やグラフなどの画像化に便利

PDF

2-2 サンプル PDF

サンプルとして配布している PDF は次のような仕様となっています。

G1

表2 サンプル PDF 判型・基本版面

判型	基本版面
新書判縦組	17行／ページ、42文字／行、基本文字サイズ：9ポイント (pt)、行送り：14.2pt

PNG



表2 サンプル PDF 判型・基本版面

判型	基本版面
新書判縦組	17行／ページ、42文字／行、基本文字サイズ：9ポイント (pt)、行送り：14.2pt



HTMLタグ付き保存

- ▶ HTMLタグ付きテキストでの保存
- ▶ bodyの下位に出力するタグは次の8種類
 - ▶ 段落タグ (p)
 - ▶ 見出しタグ (h1~h6)
 - ▶ イメージタグ (img)
 - ▶ イメージ範囲にテキストがあるとき、テキストをimgタグのalt属性値に保存できる



デモ

PDF Advanced Extractor [白馬岳.pdf]

ファイル ホーム 表示 ツール

開く 保存 手のひら スーム 前の表示に戻す ツール

貼り付け 切り取り コピー 削除 編集

自動認識 認識

見出し 保存形式 本文 段落区切り テキストを保存 画像 枠の作成


抽出順序指定 指定漏れチェック テキスト抽出

HTML保存 TEXT保存 ページ範囲

1 2 3 4 5 6 7 8 9

白馬岳

木暮理太郎



位置

今では日本アルプスの名で広く知られている飛騨山脈は、加藤理学士の説によると、凡そ南西より北十度東に向って並走せる数条の連脈から成っているものであるという。其連脈の二に白馬山脈というのがある。立山山脈との対称上から又後立山山脈とも呼ばれ、飛騨山脈中の最も長い山脈で、北は日本海岸の親不知附近から起り、越中と越後及び信濃との国境を南走して遠く飛騨国内に達しているが、中に於て越中、越後及び信濃の三国界から飛騨、信濃及び北陸の三国界附近に至る、直径にして五十五約十四里の間が主要部という可き部分であつて、最高二千九百九十米、最低二千八百八十米、平均高度は二千六百二十米に及んでゐる。そして二千八百米を超えてゐる峰は十五、六座を下らないのである。それが松本市の西の縁から大黒山を建てたやうに急に聳え立っているので、地形の相違の著しい為、二千五百米以下に於ては中山性の地帯と称す可きものに属するに拘わらず、恰も大山脈を見るが如き観を呈し、加うるに盛夏八月の候も尚お白に輝く雪田の山の嶺を飾り、雪渓が幾条となく山脈に象眼されているので、頂上附近の高山性地帯と相俟つて、一層崇高偉大なる感じを起さしめるのである。

地質

白馬山脈を構成する岩石は、大体に於て花崗岩又は之に類似した深成岩であるが、新火山岩が其間に噴出し、又古層の露出せる所も少なくない。然し調査が充分に行き届いてゐる訳ではないから、詳細に探究された訳には、新に発見する所が少なくないであらうと思ふ。白馬岳の近傍は花崗岩、蛇紋岩、古生層、玢岩及び新火山岩などで構成されてゐる。そして頂上附近は概して東側は古生層、西側は玢岩から成り、中央の一部に蛇紋岩の認められることが、地質調査所の地図に明示されてゐる。

昭和十二年二月東京地学協会発行の「白馬岳」図冊に掲載

1

レディ 位置 サイズ 233 x 268 mm 200 dpi << < ページ 1 / 9 > >> 倍率 70 % + 日



製品概要

項目	内容
製品名	PDF Advanced Extractor
バージョン	α版（仕様が変更される可能性があります）
動作環境	Windows 10 64ビット版
種別	デスクトップ製品（サーバー上で使ったり、システムに組み込むことはできません）
価格	オープン価格
発売日	未定
製品Web	https://www.antenna.co.jp/pdfae/



α版試用と評価のお願い

- ▶ 本製品のα版をお試しく下さい。
- ▶ 評価結果をお待ちしています。

次のURLから評価版をダウンロードいただけます：

<https://www.antenna.co.jp/pdfae/eval.html>

お問い合わせ先：oem@antenna.co.jp



参考資料

アンテナハウスPDF資料室

- ▶ PDFからテキストを取り出す処理について
「プログラマーから見たPDFファイル」
<https://www.antenna.co.jp/pdf/reference/pdftext.html>
- ▶ PDFテキスト抽出の諸問題
「簡単そうで簡単ではないPDFのテキスト抽出」
<https://www.antenna.co.jp/pdf/reference/text-extractor.html>

