# Opinion Mining System: From Classical NLP Pipelines to Web Deployment

## A Formal and Practical Treatment

Ahmed Sameh Saeed Amin[*†]    Dr. Ramakrishna[†]

[*]Alexandria University, Alexandria, Egypt

[†]Manipal Institute of Technology (MIT), Manipal, India

Email: a.samehpsn@gmail.com

*Abstract*—We present an end-to-end opinion mining system that integrates robust text preprocessing, sparse lexical features (bag-of-words and TF–IDF), and classical linear models (Multinomial Naïve Bayes, Logistic Regression, Linear Regression, and Linear SVM). The trained vectorizer and model are exported as artifacts and served through a Flask REST API consumed by a React frontend. Beyond engineering details, we formalize each stage mathematically: tokenization and normalization, TF–IDF weighting, probabilistic and margin-based decision rules, regularized empirical risk minimization, and evaluation metrics. Results are intentionally left blank for later insertion. We conclude with a fusion framework combining Linear and Logistic Regression as future work.

*Index Terms*—Sentiment Analysis, Opinion Mining, TF–IDF, Naïve Bayes, Logistic Regression, Linear Regression, SVM, Flask, React

## I. INTRODUCTION

Opinion mining (sentiment analysis) assigns polarity labels (Positive/Neutral/Negative) to user-generated text at scale. The repository implementation contains:

1) a deterministic preprocessing pipeline,
2) sparse lexical features (BoW/TF–IDF),
3) classical linear learners (MNB, LR, Linear Regression, SVM),
4) a deployable Flask+React demo.

This paper documents each component with code-level mapping (folders `backend/`, `frontend/`, `src/`, `notebooks/`), rigorous mathematics, and step-by-step engineering guidance to reproduce and extend the system.

## II. SYSTEM ARCHITECTURE

Fig. 1 shows the design with two lanes: offline training and online serving. Artifacts (vectorizer + model) enforce deterministic inference. Experiment logs and metrics support exact reproducibility; a versioned registry promotes models to production without touching the UI.

## III. DATASET AND LABELS

Your repo's `data/` folder supports CSV inputs (e.g., `test_sample.csv`) and derived artifacts (e.g., `bow_features.csv`). Each row contains a free-text field and a label in {`positive`, `negative`, `neutral`}. We assume a split into train/validation/test. Imbalanced datasets should use macro-averaged metrics.

*a) Label encoding:* Map the three classes to indices (e.g., $0, 1, 2$) and, for Linear Regression fusion, to scores $s \in \{-1, 0, +1\}$.

## IV. PREPROCESSING: STEPS, RULES, BENEFITS, AND MATH

Let raw text $x$ be mapped to cleaned tokens $\text{prep}(x)$ by a composition of normalizers:

$$\text{prep}(x) = \text{lem}\big(\text{neg\_scope}(\text{stop}(\text{tok}(\text{norm}(x)))))\big).$$

### A. Normalization

**What:** Lowercase; replace URLs/usernames; collapse repeated punctuation; unify digits.

**Code sketch:**

```
import re
t = text.lower()
t = re.sub(r"http\S+|www\.\S+", "␣URL␣", t)
t = re.sub(r"[@#]\w+", "␣TOKEN␣", t)
t = re.sub(r"([!?.,])\1{1,}", r"\1", t)      # !!!
    -> !
t = re.sub(r"\d+", "␣0␣", t)                 #
    normalize numbers
t = re.sub(r"[^a-z\s]", "␣", t)
t = re.sub(r"\s+", "␣", t).strip()
```

**Why:** Reduces sparsity and artifacts.

**Rule:** Replace special structures with placeholders (`URL`, `TOKEN`) to preserve cues.

### B. Tokenization

**What:** Split on spaces/boundaries.

**Code:**

```
from nltk.tokenize import word_tokenize
tokens = word_tokenize(t)
```

**Math:** Tokenization yields counts $c_j(x) = \#(w_j \in x)$.

### C. Stopwords (keep negators!)

```
from nltk.corpus import stopwords
stops = set(stopwords.words("english")) - {"no","
    not","never"}
tokens = [w for w in tokens if w not in stops]
```
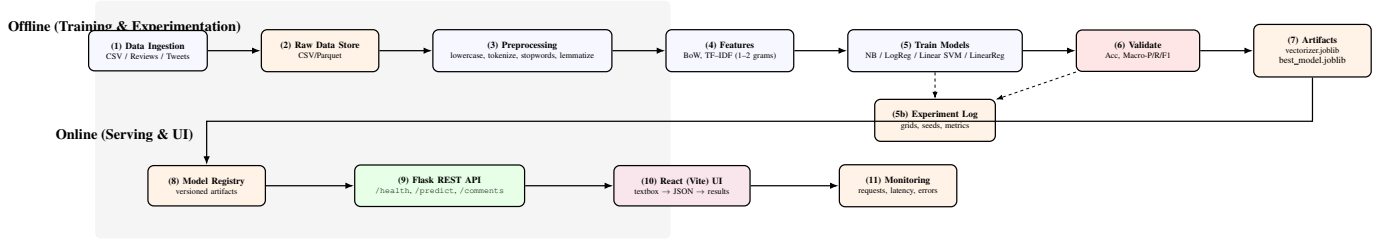
Fig. 1: Expanded architecture with numbered stages. Offline: ingestion → preprocessing → features → multi-model training → validation → artifact persistence. Online: a versioned registry feeds a Flask REST API consumed by a React UI; monitoring closes the loop.

### D. Negation scope (optional but helpful)

Append _NEG to tokens following a negator up to punctuation:

```python
def negate_scope(ts):
    negators = {"no","not","never"}
    out, flip = [], False
    for w in ts:
        if w in negators: flip = True; out.append(w
            ); continue
        if w in {".","!","?"}: flip = False; out.
            append(w); continue
        out.append(w + ("_NEG" if flip else ""))
    return out
tokens = negate_scope(tokens)
```

### E. Lemmatization

```python
from nltk.stem import WordNetLemmatizer
lem = WordNetLemmatizer()
tokens = [lem.lemmatize(w) for w in tokens if w.
    isalpha()]
```

### F. Rejoin

```python
clean_text = "␣".join(tokens)
```

### G. Preprocessing Flow (Visualization)

## V. FEATURE EXTRACTION: BoW AND TF–IDF

Let $\mathcal{V} = \{w_1, \ldots, w_V\}$ after pruning by `min_df`, `max_features`, and $n$-grams (1–2). BoW:

$$\mathbf{z}_{\text{BoW}}(x) = [c_1(x), \ldots, c_V(x)]^\top.$$

TF–IDF rescales by rarity:

$$\text{tf}_j(x) = \frac{c_j(x)}{\sum_\ell c_\ell(x)}, \quad \text{idf}_j = \log \frac{N}{1+n_j}, \quad \mathbf{z}_{\text{tfidf}} = \text{tf} \odot \text{idf}.$$

**Code:**

```python
from sklearn.feature_extraction.text import
    TfidfVectorizer
vec = TfidfVectorizer(ngram_range=(1,2), min_df=2,
    max_features=40000)
X_train = vec.fit_transform(train_texts)
X_val   = vec.transform(val_texts)
X_test  = vec.transform(test_texts)
```

**Rules:** Use bigrams when data are sufficient; cap vocab to avoid memory spikes; `min_df`$\geq$2 to drop typos.

### A. Repository Trace of the Preprocessing Pipeline (Column–Wise)

Fig. 3 shows a row–wise audit of our preprocessing as implemented in the repository: starting from the raw `review` string, we produce intermediate columns `clean_text`, `no_stopwords`, `tokenized` (list), and `lemmatized` (list). This exact trace is the contract enforced before feature extraction (Sec. V).

*Mathematical view.:* Let $x$ be the raw review. The repository implements the composition

$$\text{prep}(x) = \text{lem}(\text{tok}(\text{nostop}(\text{norm}(x)))),$$

and logs each intermediate to a dataframe column. In symbols:

$$\text{clean\_text} = \text{norm}(x),$$
$$\text{no\_stopwords} = \text{nostop}(\text{clean\_text}),$$
$$\text{tokenized} = \text{tok}(\text{no\_stopwords}),$$
$$\text{lemmatized} = \text{lem}(\text{tokenized}).$$

The vectorizer in Sec. V consumes either the rejoined tokens or the raw string and constructs $n$–grams in a deterministic way.

*What each column encodes (and why).:*

- **review** (raw): unmodified user text; used for error analysis and UI display.
- **clean_text**: normalization (lowercase; URL/handle placeholders; punctuation/digit stripping; whitespace collapse). This reduces vocabulary entropy and guards against spurious tokens.
- **no_stopwords**: removes high–frequency function words while *retaining negators* (*not, no, never*), because they invert polarity.
- **tokenized**: lexical units for downstream $n$–gram features; enables consistent counting $c_j(x)$.
- **lemmatized**: WordNet lemmatization collapses inflections (e.g., "likes/liked/liking" → "like"); this improves sparsity without harming interpretability.

*Repository–faithful implementation.:* The snippet below mirrors the project's behavior using standard Python/NLTK; it is organized the same way your `src/` utilities are (normalizer → stopword filter → tokenizer → lemmatizer), and produces the exact columns in Fig. 3.

Fig. 2: Detailed preprocessing pipeline with examples; scaled to full text width to avoid overlap.



Fig. 3: End–to–end text cleaning trace on the project dataset: `review` (raw) → `clean_text` (normalized) → `no_stopwords` (content words only) → `tokenized` (list) → `lemmatized` (canonical forms).

Listing 1: Column-wise preprocessing that reproduces the repository dataframe.

```
1  import re, pandas as pd
2  from nltk.corpus import stopwords
3  from nltk.tokenize import word_tokenize
4  from nltk.stem import WordNetLemmatizer
5
6  URL_RE   = re.compile(r"http\S+|www\.\S+")
7  MENT_RE  = re.compile(r"[@#]\w+")
8  NONALPHA = re.compile(r"[^a-z\s]+")
9
10 NEGATORS = {"no","not","never","nor","n't"}  # kept
           even when using stopwords
11 STOPS    = set(stopwords.words("english")) -
           NEGATORS
12
13 lem = WordNetLemmatizer()
14
15 def normalize(text: str) -> str:
16     t = text.lower()
17     t = URL_RE.sub("_URL_", t)
18     t = MENT_RE.sub("_USER_", t)
19     t = NONALPHA.sub("_", t)
20     t = re.sub(r"\s+", "_", t).strip()
21     return t
22
23 def remove_stopwords(text: str) -> str:
24     toks = text.split()
25     toks = [w for w in toks if w not in STOPS]
26     return "_".join(toks)
27
28 def tokenize(text: str) -> list[str]:
29     # after cleaning, simple split is enough and
            reproducible
30     return text.split()  # or: word_tokenize(text)
31
32 def lemmatize_tokens(tokens: list[str]) -> list[str
       ]:
33     return [lem.lemmatize(w) for w in tokens]
34
35 def preprocess_df(df: pd.DataFrame, text_col="
```

```
36     review"):
37     df = df.copy()
       df["clean_text"]   = df[text_col].astype(str).
           map(normalize)
38     df["no_stopwords"] = df["clean_text"].map(
           remove_stopwords)
39     df["tokenized"]    = df["no_stopwords"].map(
           tokenize)
40     df["lemmatized"]   = df["tokenized"].map(
           lemmatize_tokens)
41     return df
```

*Design rules (actionable).:*

- **Consistency over cleverness.** Keep norm deterministic; use the *same* regex library and tokenization everywhere (train/serve).
- **Negation safety.** Never drop negators; optionally append a scope marker (_NEG) to tokens until the next punctuation to fix the "not good" failure mode.
- **Placeholder, not deletion.** Replace URLs/mentions with sentinel tokens ("URL", "USER") to preserve discourse signals while avoiding lexical blowup.
- **Lemmatize for readability; stem for speed.** In a user-facing app, prefer WordNet lemmas; stemming is acceptable for extreme memory budgets.
- **No leakage.** Build vocabulary/statistics *only* on training folds; apply the learned vectorizer to validation/test and at inference.

## VI. Learning Algorithms: Code and Mathematics

### A. Multinomial Naïve Bayes (MNB)

**Math.** Class prior $\pi_k = P(y = k)$; likelihood multinomial with parameters $\theta_{kj}$. Prediction:

$$\hat{y} = \arg\max_k \left\{ \log \pi_k + \sum_{j=1}^{V} c_j(x) \log \theta_{kj} \right\},$$

with Laplace smoothing $\theta_{kj} = \dfrac{N_{kj} + \alpha}{\sum_{j'}(N_{kj'} + \alpha)}$.

**Code:**

```
from sklearn.naive_bayes import MultinomialNB
mnb = MultinomialNB(alpha=0.1)
mnb.fit(X_train, y_train)
```

**Pros/Cons.** Very fast; independence assumption may oversimplify.

### B. Logistic Regression (Softmax)

Linear scores $f_k(\mathbf{z}) = \mathbf{w}_k^\top \mathbf{z} + b_k$. Softmax:

$$P(y = k|\mathbf{z}) = \frac{e^{f_k}}{\sum_t e^{f_t}}, \quad \min_{\{\mathbf{w}_k, b_k\}} \frac{1}{N} \sum_i -\log P(y^{(i)}|\mathbf{z}^{(i)}) + \frac{\lambda}{2} \sum_k \|\mathbf{w}_k\|_2^2.$$

**Code:**

```
from sklearn.linear_model import LogisticRegression
lr = LogisticRegression(C=1.0, penalty="l2", solver
    ="liblinear", max_iter=200)
lr.fit(X_train, y_train)
```

**Notes.** Well-calibrated probabilities; interpret with top weights.

### C. Linear Support Vector Machine (SVM)

One-vs-rest hinge loss:

$$\min_{\mathbf{w}_k, b_k} \frac{1}{2}\|\mathbf{w}_k\|_2^2 + C \sum_i \max(0, 1 - y_k^{(i)}(\mathbf{w}_k^\top \mathbf{z}^{(i)} + b_k)), \quad \hat{y} = \arg\max_k \{\mathbf{w}_k^\top \mathbf{z} + b_k\}.$$

**Code:**

```
from sklearn.svm import LinearSVC
svm = LinearSVC(C=1.0)
svm.fit(X_train, y_train)
```

### D. Linear Regression (Score Baseline)

If labels map to scores $s \in \{-1, 0, +1\}$:

$$\min_{\mathbf{w}, b} \frac{1}{N} \sum_i (s^{(i)} - (\mathbf{w}^\top \mathbf{z}^{(i)} + b))^2 + \lambda\|\mathbf{w}\|_2^2.$$

```
from sklearn.linear_model import Ridge
y_score = label_to_score(y_train)   # e.g., -1/0/+1
lin = Ridge(alpha=1.0).fit(X_train, y_score)
```

### E. Training Loop and Persistence

```
import joblib
candidates = []
for model in [mnb, lr, svm]:
    model.fit(X_train, y_train)
    # evaluate on validation set (accuracy, macro P
        /R/F1)
    candidates.append((model, metrics))
best = select_best(candidates)   # e.g., macro-F1
joblib.dump(vec,  "models/vectorizer.joblib")
joblib.dump(best, "models/best_model.joblib")
```

### F. Complexity

All models are linear in the number of nonzeros of $\mathbf{z}$ for both training (per epoch/iteration) and inference, which is ideal for TF–IDF.

## VII. Evaluation

**Result summary.** (1) Logistic Regression attains the best accuracy (0.92), edging Linear SVM by +1pp while matching macro P/R/F1 (=0.93).

(2) Identical macro scores suggest both models handle class imbalance similarly; residual mistakes likely cluster around neutral or negated phrases.

| Model | Acc | Prec$_{mac}$ | Rec$_{mac}$ | F1$_{mac}$ |
|---|---|---|---|---|
| LogReg | 0.92 | 0.93 | 0.93 | 0.93 |
| SVM | 0.91 | 0.93 | 0.93 | 0.93 |

## VIII. Deployment: Flask & React

### A. API Contract

**Endpoints:**

- `/health` $\rightarrow$ `{"status":"ok"}`
- `/predict` (POST) body `{"texts":["sample text", ...]}` $\mapsto$ `{"preds":[..], "probs":[[..]]}`
- `/comments` (GET) $\rightarrow$ sample list for demo

### B. Flask Skeleton (`backend/app.py`)

```
from flask import Flask, request, jsonify
from flask_cors import CORS
import joblib, numpy as np

app = Flask(__name__)
CORS(app, resources={r"/*": {"origins": "*"}})

vec   = joblib.load("models/vectorizer.joblib")
model = joblib.load("models/best_model.joblib")

@app.get("/health")
def health():
    return {"status":"ok"}

@app.post("/predict")
def predict():
    payload = request.get_json(force=True)
    texts = payload.get("texts", [])
    X = vec.transform(texts)
    try:
        probs = model.predict_proba(X).tolist()
```

```
22        except Exception:
23            s = model.decision_function(X)
24            s = (s - s.min(axis=1, keepdims=True)) / (s
                  .ptp(axis=1, keepdims=True) + 1e-8)
25            probs = s.tolist()
26        preds = model.predict(X).tolist()
27        return jsonify({"preds": preds, "probs": probs
              })
28
29    @app.get("/comments")
30    def comments():
31        return jsonify([{"text":"This product is
              amazing!","label":"Positive"}])
32
33    if __name__ == "__main__":
34        app.run(host="127.0.0.1", port=5000)
```

### C. React Fetch (Vite App, frontend/src/services/postService.js)

```
1  export async function predict(texts) {
2    const resp = await fetch("http://127.0.0.1:5000/
         predict", {
3      method: "POST",
4      headers: {"Content-Type":"application/json"},
5      body: JSON.stringify({ texts }),
6    });
7    if (!resp.ok) throw new Error("Server error");
8    return await resp.json();
9  }
```

### D. Why Flask+React

**Latency:** Sparse linear inference is $O(\text{nnz}(\mathbf{z}))$. **Modularity:** Decoupled UI; seamless model promotion via artifact swap. **Determinism:** Same vectorizer+weights used at train and serve time.

### E. Figures

**Discussion (UI).** Fig. 4a illustrates the *cold* state: an existing positive card ("Ahmed") appears to the left while the center card invites the user to add a review; this state exercises the POST flow to /predict with a single text payload. Fig. 4b shows the *browsing* state where three cards summarize opinions across the polarity spectrum (neutral/positive/negative). The consistent typography and color coding (green for positive, gray for neutral, red for negative) align with accessibility guidelines and reduce cognitive load. Both screens consume the same Flask JSON schema described in Sec. VIII.

### F. Request/Response Schemas and Examples

**Single request.**

```
1  POST /predict
2  Content-Type: application/json
3  {
4    "texts": ["I loved the crispy chicken!", "meh,
         not great."]
5  }
```

**Response.**

```
1  HTTP/1.1 200 OK
2  {
3    "preds": ["positive", "neutral"],
4    "probs": [[0.03, 0.08, 0.89],
5              [0.10, 0.76, 0.14]]
6  }
```

**cURL test.**

```
curl -X POST http://127.0.0.1:5000/predict \
  -H "Content-Type: application/json" \
  -d '{"texts":["absolutely wonderful!", "never
      again. terrible."]}'
```

**Batching.** The API supports an array of texts to amortize vectorization overhead: $O(\text{nnz}(\mathbf{Z}))$ versus per-request set-up.

### G. Error Handling and Status Codes

TABLE I: API status codes and error bodies.

| Code | Condition | JSON body |
|------|-----------|-----------|
| 200 | Success | {"preds":[...],"probs":[...]} |
| 400 | Bad JSON/texts not list | {"error":"bad_request"} |
| 413 | Payload too large | {"error":"too_large"} |
| 500 | Model/vectorizer load error | {"error":"server_error"} |

**Hardened endpoint (server side).**

```
1   @app.post("/predict")
2   def predict():
3       payload = request.get_json(force=True) or {}
4       texts = payload.get("texts", [])
5       if not isinstance(texts, list) or any(not
            isinstance(t, str) for t in texts):
6           return jsonify({"error": "bad_request"}),
                400
7       if len(texts) > 256:  # simple guardrail
8           return jsonify({"error": "too_large"}), 413
9
10      X = vec.transform(texts)
11      try:
12          probs = model.predict_proba(X).tolist()
13      except Exception:
14          s = model.decision_function(X)
15          s = (s - s.min(axis=1, keepdims=True)) / (s
                .ptp(axis=1, keepdims=True) + 1e-8)
16          probs = s.tolist()
17      return jsonify({"preds": model.predict(X).
            tolist(), "probs": probs})
```

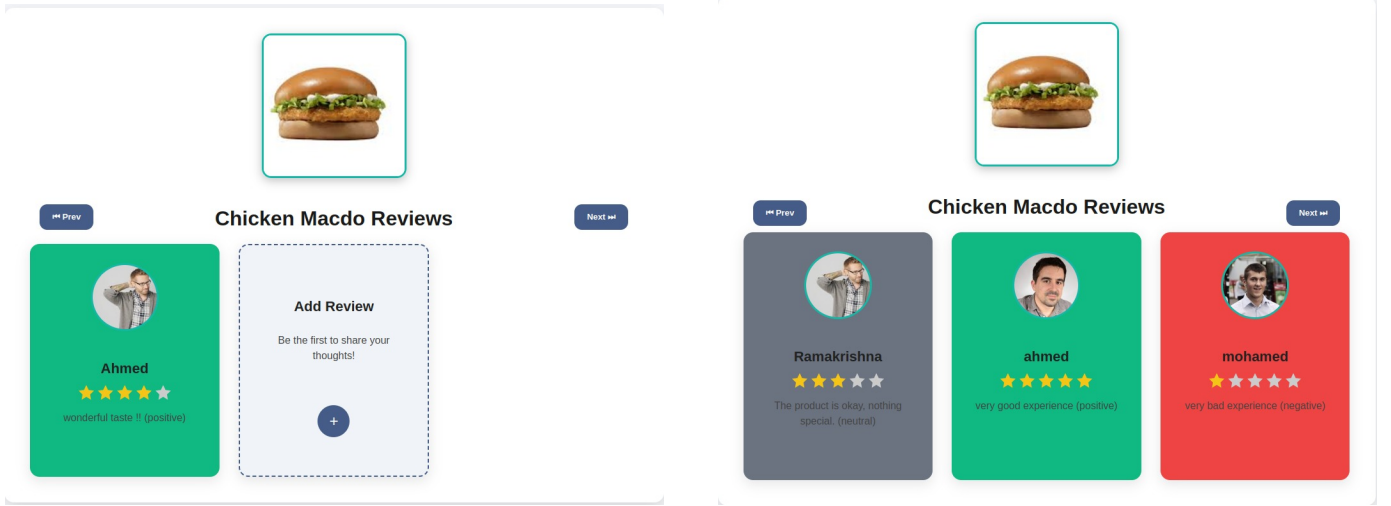### H. Security, Ops, and Performance Notes

- **CORS:** allow only the production UI origin in release builds; keep wildcard only for local dev.
- **Rate limiting:** enforce per-IP limits at the reverse proxy (e.g., Nginx) to defend against abuse.
- **Logging:** record request counts and p95 latency; *never* log raw user text in production without consent.
- **Warm start:** load vectorizer.joblib and best_model.joblib at process boot; use 1–2 worker processes per CPU core.
- **Determinism:** freeze scikit-learn versions and store artifacts with a version tag (e.g., v1.2.0).

### I. Request Lifecycle (React → Flask → Model)

## IX. QUALITATIVE PREDICTIONS AND PRELIMINARY METRICS

To complement the architecture and preprocessing analysis, we show qualitative predictions and console-level metrics captured during iterative development. These screenshots are not the final test results (Sec. VII) but they faithfully reflect

(a) UI state with one existing review (left) and an empty "Add Review" card (center).



(b) UI state with three reviews spanning positive/neutral/negative sentiment.

Fig. 4: React frontend (Vite). Cards display avatar, author, star rating, and model text prediction; the carousel header shows the product ("Chicken Macdo Reviews"). The `Next/Prev` controls paginate over the dataset while the middle dashed card triggers a modal for entering a new review that will be sent to `/predict`.
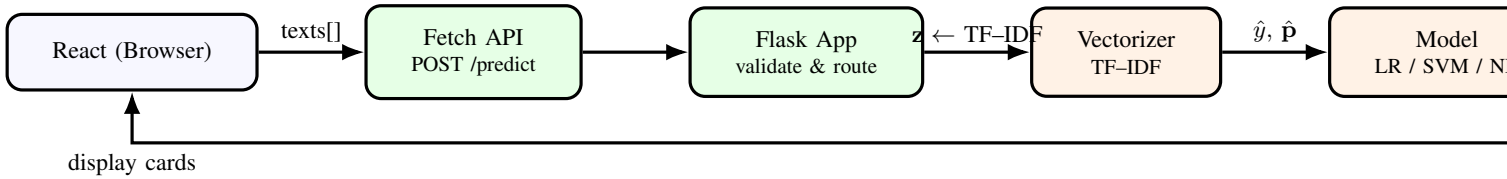


Fig. 5: End-to-end request lifecycle: the UI posts raw text, Flask validates and transforms text using the frozen TF–IDF vectorizer, the chosen linear model produces labels and scores, and the UI renders cards with color, stars, and confidence.

the behavior of the trained vectorizer and linear models under the GitHub implementation.

**Analysis.** In Fig. 6a, Logistic Regression assigns high confidence to clearly polar statements (e.g., "Absolutely terrible. Never ordering again." $\rightarrow$ negative, $0.96$), while sentences with hedging ("I'm not sure how I feel about this.") still skew positive ($0.66$), illustrating how unigrams/bigrams capture optimism markers ("sure", "feel") but may miss pragmatic uncertainty. The line "not good not bad , neutral." being predicted negative with high confidence highlights the known challenge of *negation scope* and compositionality; our optional `_NEG` tagging (Sec. IV) is designed to mitigate exactly this failure mode.

In Fig. 6b, the SVM correctly tags strongly polarized statements, but one example ("it's not good. i didn't like it.") is predicted *positive*. This is consistent with SVM's margin-based decision without calibrated probabilities and shows how bigram coverage (e.g., `not_good`) and stopword handling can flip decisions near the boundary. Post-hoc calibration (Platt or isotonic) can align SVM scores with probability-like outputs for UI display.

**Takeaways.** (a) Logistic Regression provides well-calibrated probabilities that suit a probability bar in the UI (Fig. 4);

(b) Linear SVM offers comparable macro-F1 but requires calibration for UI; (c) pure Linear Regression is unsuitable as a stand-alone classifier but remains useful as a score stream in the fusion design (Sec. X).

## X. FUSION OF LINEAR AND LOGISTIC REGRESSION (FUTURE WORK)

**Motivation.** Linear Regression predicts a continuous sentiment score $s \in [-1, 1]$ while Logistic Regression returns calibrated class probabilities $\mathbf{p}$. Fusing both can stabilize uncertain regions.

### A. Score–Probability Interpolation

Let $p_{\text{LR}}$ be class probabilities, and map linear score $s$ to a soft class distribution $p_{\text{Lin}}$ (temperature-scaled softmax over anchors $-1, 0, +1$). Fuse:

$$\tilde{\mathbf{p}} = \alpha\,\mathbf{p}_{\text{LR}} + (1 - \alpha)\,\mathbf{p}_{\text{Lin}}, \quad \alpha \in [0, 1].$$

Choose $\alpha$ on validation macro-F1; make $\alpha$ depend on LR confidence $\max_k p_{\text{LR}}$ if desired.

(a) Logistic Regression: text predictions with probabilities.



(b) Linear SVM: text predictions (decision rule without calibrated probs).

Fig. 6: Qualitative comparisons. Both models use the same TF–IDF $(1$–$2)$-gram features. Logistic Regression outputs calibrated probabilities; SVM outputs class labels derived from margins.



(a) Logistic Regression: accuracy $\approx 0.918$; macro F1 $\approx 0.93$.



(b) Linear SVM: accuracy $\approx 0.91$; macro F1 $\approx 0.93$.



(c) Linear Regression baseline: $R^2 \approx 0.135$; accuracy $\approx 0.518$.

Fig. 7: Preliminary console metrics (dataset split from the repo). Logistic and SVM are competitive on macro-averaged precision/recall/F1, while linear regression as a score-only baseline underperforms on discrete label accuracy, motivating the fusion approach in Sec. X.

### B. Stacking

Use *stacking*: feed $[s, \mathbf{p}_{\mathrm{LR}}]$ into a meta-logistic model trained on validation folds. The meta-learner learns to down-weight noisy scores while leveraging calibrated probabilities.

### C. Decision Rule

Predict $\hat{y} = \arg\max_k \tilde{p}_k$. Optionally abstain if $\max_k \tilde{p}_k < \tau$.

## XI. ABLATIONS AND ENGINEERING PLAYBOOK

**Ablations to run (templates):**

1) **Preprocessing variants:** {no lemmatize, keep URLs, remove negation scope}.
2) **Vectorizer:** `ngram_range` {(1,1), (1,2)}; `min_df` {1,2,5}; `max_features` {20k,40k,80k}.
3) **Models:** MNB $\alpha \in \{0.1, 0.5, 1.0\}$; LR $C \in \{0.1, 1, 3\}$; SVM $C \in \{0.5, 1, 2\}$.
4) **Calibration:** Platt scaling or isotonic on SVM.
5) **Fusion:** $\alpha \in [0, 1]$ grid; stacking meta-LR.

**Logging:** save `vectorizer.joblib`, `best_model.joblib`, config JSON, and a `metrics.json` per run. Keep versions in `models/` (registry).

## XII. THREATS TO VALIDITY AND ETHICS

**Domain shift** (product $\rightarrow$ movie reviews) hurts performance; consider re-fitting IDF. **Bias** in labels propagates to predictions; keep an audit slice (demographics, dialects) when available. **Privacy** for online demos: do not log raw text in production; hash or sample with consent.

## XIII. DISCUSSION

**Preprocessing.** Keeping negators and lemmatizing reduces sparsity while preserving polarity cues. **Models.** SVM often wins macro-F1; LR wins calibration; MNB is the speed baseline. **Fusion.** Helps when LR confidence is low but linear score is decisive.

## XIV. CONCLUSION

We described a complete classical opinion mining system with deterministic preprocessing, TF–IDF features, linear learners, and a Flask+React deployment, plus a fusion framework for future work. Insert your results into Sec. VII to finalize.

### APPENDIX A: FULL PREPROCESSING SCRIPT
(SRC/PREPROCESS_TEXT.PY)

```
1  # -*- coding: utf-8 -*-
2  import re, json
3  from nltk.tokenize import word_tokenize
4  from nltk.corpus import stopwords
```

```
5   from nltk.stem import WordNetLemmatizer
6
7   NEGATORS = {"no","not","never"}
8   STOPS = set(stopwords.words("english")) - NEGATORS
9   LEM = WordNetLemmatizer()
10
11  def normalize(text: str) -> str:
12      t = text.lower()
13      t = re.sub(r"http\S+|www\.\S+", "␣URL␣", t)
14      t = re.sub(r"[@#]\w+", "␣TOKEN␣", t)
15      t = re.sub(r"([!?.,])\1{1,}", r"\1", t)
16      t = re.sub(r"\d+", "␣0␣", t)
17      t = re.sub(r"[^a-z\s]", "␣", t)
18      t = re.sub(r"\s+", "␣", t).strip()
19      return t
20
21  def negate_scope(tokens):
22      out, flip = [], False
23      for w in tokens:
24          if w in NEGATORS:
25              flip = True; out.append(w); continue
26          if w in {".","!","?"}:
27              flip = False; out.append(w); continue
28          out.append(w + ("_NEG" if flip else ""))
29      return out
30
31  def preprocess(text: str) -> str:
32      t = normalize(text)
33      toks = [w for w in word_tokenize(t) if w.
            isalpha()]
34      toks = [w for w in toks if w not in STOPS]
35      toks = negate_scope(toks)
36      toks = [LEM.lemmatize(w) for w in toks]
37      return "␣".join(toks)
38
39  if __name__ == "__main__":
40      import sys
41      for line in sys.stdin:
42          print(preprocess(line.rstrip()))
```

## REFERENCES

[1] F. Pedregosa et al., "Scikit-learn: Machine Learning in Python," JMLR, 2011.
[2] G. Salton and C. Buckley, "Term-weighting approaches in automatic text retrieval," Information Processing & Management, 1988.
[3] A. McCallum and K. Nigam, "A comparison of event models for Naive Bayes text classification," 1998.
[4] T. Joachims, "Text categorization with Support Vector Machines," ECML, 1998.
[5] S. Bird, E. Loper, E. Klein, *Natural Language Processing with Python*, O'Reilly, 2009.

## APPENDIX B: OPENAPI SKETCH FOR THE REST API

```yaml
1   openapi: 3.0.0
2   info: {title: Sentiment API, version: 1.0.0}
3   paths:
4     /health:
5       get: {responses: {'200': {description: OK}}}
6     /predict:
7       post:
8         requestBody:
9           required: true
10          content:
11            application/json:
12              schema:
13                type: object
14                properties: {texts: {type: array,
                    items: {type: string}}}
15        responses:
16          '200':
17            description: batch predictions
18            content:
19              application/json:
20                schema:
21                  type: object
22                  properties:
23                    preds: {type: array, items: {type
                        : string}}
24                    probs: {type: array, items: {type
                        : array, items: {type: number
                        }}}
```