

- Useful for 1000s of authors/publishers where one-hot would explode
- **Critical:** Must use out-of-fold encoding to avoid leakage (don't use a sample's own price in its encoding)
- Add smoothing for rare categories (blend category mean with global mean)

Embedding layers (for neural networks):

python

- `# Keras example`
- `author_input = Input(shape=(1,))`
- `author_embedding = Embedding(input_dim=num_authors,`
`output_dim=10)(author_input)`

Now each author is a learned 10-dim vector

- Network learns that similar authors should have similar embeddings
- Captures latent structure (e.g., "literary fiction authors" cluster together)

Text Feature Engineering

Bag-of-Words / TF-IDF:

python

- `from sklearn.feature_extraction.text import TfidfVectorizer`
- `tfidf = TfidfVectorizer(max_features=500, ngram_range=(1,2),`
`stop_words='english')`

`X_text = tfidf.fit_transform(df['description'])`

- Converts text to sparse numeric matrix
- TF-IDF weighs words by importance (rare words get higher weight than common ones)
- N-grams capture phrases ("first edition" as bigram more meaningful than separate words)

Keyword extraction:

python

- `df['has_signed'] = df['title'].str.contains('signed', case=False).astype(int)`
- `df['has_first_edition'] = df['title'].str.contains('first edition', case=False).astype(int)`

`df['has_damage_keywords'] = df['description'].str.contains('water|torn|stain',`
`case=False).astype(int)`

- Simple but effective: presence of specific keywords as binary features
- Domain knowledge drives which keywords to extract

Topic modeling:

python

- `from sklearn.decomposition import LatentDirichletAllocation`
- `lda = LatentDirichletAllocation(n_components=20)`
- `topic_features = lda.fit_transform(bow_matrix)`

Each book now has a 20-dim vector of topic proportions

- Learns latent themes in text (e.g., Topic 5 = "romance, love, relationship")
- Can reveal genre/subject without explicit labels

Pre-trained embeddings:

python

- `# Using sentence transformers`
- `from sentence_transformers import SentenceTransformer`
- `model = SentenceTransformer('all-MiniLM-L6-v2')`
- `synopsis_embeddings = model.encode(df['synopsis'].tolist())`

Each synopsis → 384-dim embedding

- Captures semantic meaning
- "Epic fantasy adventure" and "sword and sorcery quest" will have similar embeddings
- Can feed these embeddings into XGBoost or use as NN input

Sentiment analysis:

python

- `from textblob import TextBlob`
- `df['review_sentiment'] = df['reviews'].apply(lambda x: TextBlob(x).sentiment.polarity)`

Polarity ranges from -1 (negative) to +1 (positive)

- Positive reviews may correlate with higher prices (demand signal)
- Can be noisy; aggregate across multiple reviews more reliable

Interaction Features

Manual interactions (for linear models):

python

- `df['condition_x_age'] = df['condition_ordinal'] * df['book_age']`
- `df['price_per_page'] = df['list_price'] / df['num_pages']`

```
df['rating_x_num_ratings'] = df['avg_rating'] * np.log1p(df['num_ratings'])
```

- Captures that effect of one variable depends on another
- E.g., condition matters more for expensive books
- Tree models learn these automatically, but specifying helps linear models

Cross-features (for specific domains):

python

- `# For textbooks specifically`
- `df['is_latest_edition'] = (current_year - df['publication_year'] <= 2).astype(int)`

```
df['textbook_premium'] = df['is_textbook'] * df['is_latest_edition']
```

Handling Missing Data

Imputation strategies:

python

- `# For numeric: fill with median or -1 (if -1 is impossible, signals missingness)`
- `df['num_ratings'].fillna(df['num_ratings'].median(), inplace=True)`
-
- `# For categorical: create explicit "Unknown" category`
- `df['genre'].fillna('Unknown', inplace=True)`
-
- `# Add missing indicator as separate feature`
- `df['rating_missing'] = df['avg_rating'].isnull().astype(int)`

```
df['avg_rating'].fillna(df['avg_rating'].mean(), inplace=True)
```

Why missing indicators help:

- Missingness itself may be informative (no ratings = obscure book = lower price)
- Lets model learn different behavior for imputed vs. real values

3.7 Data Quality and Leakage Prevention

Common Data Quality Issues

1. **Duplicates:** Same book listed multiple times at different prices
 - **Solution:** Deduplicate or include all as separate training samples (reflects price variance)
2. **Outliers:** \$0.01 listings (shipping scams) or \$10,000 (pricing errors)
 - **Solution:** Cap prices at reasonable bounds (e.g., 0.1 percentile to 99.9 percentile) or remove if clearly errors
3. **Inconsistent condition grading:** Different sellers rate "Good" differently
 - **Challenge:** Hard to standardize across sources
 - **Mitigation:** Train separate models per platform, or include seller rating as feature (trusted sellers more reliable)
4. **Temporal drift:** Market conditions change (COVID spike in home library demand)
 - **Solution:** Weight recent data more heavily, or train on recent window only
5. **Selection bias:** Only seeing sold items (successful listings)
 - Unsold items had asking price too high (not observed in training)
 - **Implication:** Model may overestimate prices
 - **Mitigation:** If available, include unsold listings with a different target or build separate "will it sell at price X?" model

Feature Leakage Prevention

What is leakage?

Including information in features that wouldn't be available at prediction time, or that directly encodes the target.

Examples of leakage:

✗ Number of bids (for auction price prediction at listing time)

- Only known *after* auction runs
- Alternative: Predict number of bids separately, or don't use

✗ Final sale indicator or "sold" flag

- Directly reveals if price was acceptable
- Alternative: Remove from features

✗ Target encoding using full dataset (including test samples)

- Uses information from test set to create training features
- Alternative: Use only training fold data for encoding

✗ Historical price using future sales

- Computing "average sale price for this book" including sales that happened *after* the listing you're predicting
- Alternative: For each sample, only use sales that occurred *before* that timestamp

X Shipping tracking number, transaction ID

- Only exists after sale
- Alternative: Remove

Validation of leak-free features: Ask for each feature: "Could I obtain this value at the moment I need to make a prediction (when creating a listing)?"

- Book metadata (title, author, year): Yes
- Current Amazon lowest price: Yes (if looking it up in real-time)
- Number of current competing offers: Yes
- This book's past average sale price (from 30+ days ago): Yes (from historical database)
- Buyer's username for this sale: X No (only known after purchase)

Cross-validation must respect time: If data has timestamps:

```
python
    o # BAD: Random split (future data can leak into past predictions)
    o X_train, X_test = train_test_split(data, test_size=0.2, random_state=42)
    o
    o # GOOD: Time-based split
    o split_date = '2024-01-01'
    o train = data[data['sale_date'] < split_date]

test = data[data['sale_date'] >= split_date]
```

For time series, use **forward chaining cross-validation**:

- Fold 1: Train on Jan-Mar, validate on Apr
- Fold 2: Train on Jan-Jun, validate on Jul
- Fold 3: Train on Jan-Sep, validate on Oct
- etc.

This simulates real-world deployment where you train on past data and predict future.

4. Model Training, Evaluation, and Deployment

4.1 Evaluation Metrics: Beyond Accuracy

For price prediction, choosing the right evaluation metric is crucial—it defines what "good" means and guides optimization.

Mean Absolute Error (MAE)

python

$$\text{MAE} = (\frac{1}{n}) * \sum |\text{predicted} - \text{actual}|$$

Interpretation: Average dollar error

- E.g., MAE = \$2.50 means predictions are off by \$2.50 on average

Strengths:

- **Intuitive:** Directly interpretable in currency units
- **Robust to outliers:** Each error contributes linearly (one huge error doesn't dominate)
- **Equal weight:** Under-predictions and over-predictions penalized equally

Limitations:

- Doesn't distinguish between \$2 error on a \$5 book (40% error) vs. \$2 error on \$50 book (4% error)
- May not reflect business impact (underpricing loses revenue, overpricing loses sales)

When to use: When absolute dollar accuracy matters regardless of price magnitude.

Root Mean Squared Error (RMSE)

python

$$\text{RMSE} = \sqrt{(\frac{1}{n}) * \sum (\text{predicted} - \text{actual})^2}$$

Interpretation: Square root of average squared error (same units as price)

Strengths:

- **Penalizes large errors** more heavily due to squaring
- Sensitive to outliers (useful if big misses are especially costly)
- Common in regression (directly optimized by many algorithms)

Limitations:

- **Less interpretable** than MAE (what does RMSE = \$5 mean intuitively?)
- Can be dominated by a few large errors

- Still doesn't account for relative error

When to use: When large errors are disproportionately harmful (e.g., grossly mispricing rare collectibles)

Root Mean Squared Log Error (RMSLE)

python

$$\text{RMSLE} = \sqrt{\frac{1}{n} \sum (\log(\text{predicted} + 1) - \log(\text{actual} + 1))^2}$$

Interpretation: RMSE in log-space (measures relative/percentage error)

Strengths:

- **Penalizes relative error:** 100% error (predicting \$10 for \$5 book) matters similarly whether book is cheap or expensive
- **Reduces outlier impact:** Log compression means \$1000 error on \$5000 book matters less than in RMSE
- **Asymmetric:** Penalizes under-prediction more than over-prediction
 - Predicting \$5 when true is \$10: $\log(6) - \log(11) \approx -0.606$
 - Predicting \$15 when true is \$10: $\log(16) - \log(11) \approx 0.376$
 - Under-prediction error is $\sim 1.6 \times$ larger (economically sensible: underpricing loses revenue)

Limitations:

- Less intuitive (what does RMSLE = 0.25 mean?)
- Requires all prices > 0 (hence the +1)

Conversion to percentage:

- $\text{RMSLE} \approx 0.20 \rightarrow$ roughly 22% average relative error
- $\text{RMSLE} \approx 0.40 \rightarrow$ roughly 49% average relative error

When to use:

- **Skewed price distributions** (books range \$1 to \$500)
- When **percentage error matters more than absolute** (better to be $\pm 10\%$ wrong than $\pm \$5$ wrong)
- Marketplace pricing (Mercari, eBay) where relative accuracy is key

Why Mercari used it: "The use of logarithmic transformation helps reduce the impact of extreme values on the error metric, making it more robust to outliers."

Mean Absolute Percentage Error (MAPE)

python

$$\text{MAPE} = (\text{100}/n) * \sum |(\text{predicted} - \text{actual}) / \text{actual}|$$

Interpretation: Average percentage error

Strengths:

- **Highly interpretable:** "Predictions are off by 8% on average"
- **Scale-invariant:** Can compare across different markets or products

Limitations:

- **Undefined for actual = 0** (division by zero)
- **Asymmetric:** Heavily penalizes under-prediction
 - Predicting \$5 for \$10 book: $|5-10|/10 = 50\%$ error
 - Predicting \$15 for \$10 book: $|15-10|/10 = 50\%$ error (same)
 - But \$15 prediction is further from truth in absolute terms
- **Sensitive to low prices:** \$1 error on \$2 book is 50% error, while \$10 error on \$100 book is 10%

When to use: When dealing with prices in similar ranges (no near-zero values) and percentage error is the key business metric.

R² (Coefficient of Determination)

python

- $R^2 = 1 - (SS_{\text{residual}} / SS_{\text{total}})$
- $SS_{\text{residual}} = \sum (\text{actual} - \text{predicted})^2$

$$SS_{\text{total}} = \sum (\text{actual} - \text{mean(actual)})^2$$

Interpretation: Proportion of variance explained

- $R^2 = 1$: Perfect predictions
- $R^2 = 0$

: Model no better than always predicting the mean

- $R^2 < 0$: Model worse than mean (possible on test set)

Strengths:

- **Standardized metric** (0-1 scale) allows comparison across datasets
- **Statistical interpretation:** What fraction of price variability does the model capture?
- Common in academic research (facilitates literature comparison)

Limitations:

- **Not directly interpretable** for business decisions ("We explain 65% of variance" doesn't tell sellers what price to list)
- Can be **misleading with skewed distributions** (high R^2 doesn't guarantee good predictions on rare books)
- **Doesn't indicate direction** of errors (could systematically over/under-predict)
- Can be **inflated by outliers** in the dependent variable

Empirical benchmarks from literature:

- Linear regression with basic features: $R^2 = 0.13-0.19$
- Random Forest with rich features: $R^2 = 0.40-0.45$
- Gradient Boosting with text features: $R^2 = 0.45-0.65$
- Deep learning multimodal: $R^2 = 0.50-0.70$ (in some domains)

When to use:

- **Academic research** (standard reporting metric)
- **Model comparison** (which features/algorithms improve explanatory power?)
- **Diagnostic tool** (low R^2 indicates missing important features)

Business Metrics: Beyond Statistical Accuracy

For a production pricing system, statistical metrics are necessary but not sufficient. Consider:

1. Calibration (Uncertainty Quantification)

- **Question:** Do predicted confidence intervals contain true prices at the expected rate?
- **Metric:** If we predict "80% confident price is \$20-30", do 80% of actuals fall in that range?
- **Why it matters:** Users need to know when the model is uncertain (rare collectibles vs. common textbooks)

Test:

```
python
    # For 80% prediction intervals
    predictions_with_intervals = model.predict_interval(X_test, confidence=0.80)
    coverage = np.mean((y_test >= predictions_with_intervals[:, 0]) &
                      (y_test <= predictions_with_intervals[:, 1]))
print(f"80% interval coverage: {coverage:.2%}") # Should be ~80%
```

2. Sale Success Rate

- **Question:** What percentage of items listed at suggested price actually sell within 30 days?
- **Why it matters:** Accurate price prediction is useless if items don't sell (could be overpricing)
- **A/B test:** Compare listings using model suggestions vs. seller's intuition

3. Revenue Impact

- **Question:** Does the model increase total seller revenue vs. baseline?
- **Trade-off:** Could suggest higher prices (more revenue per sale) but lower velocity (fewer sales)
- **Metric:** Total revenue = $\Sigma(\text{sale_price} \times \text{sale_indicator})$ across all listings
- **Optimal strategy:** May not be the most "accurate" but the most profitable

4. Time-to-Sale

- **Question:** How long do items take to sell at suggested prices?
- **Trade-off:** Lower price → faster sale vs. Higher price → wait for right buyer
- **User preference:** Some sellers want quick cash, others willing to wait

5. Stratified Performance

- **Question:** Does model work well across all segments, or only for popular books?
- **Analysis:** Report MAE/RMSE separately for:
 - Textbooks vs. Fiction vs. Collectibles
 - Different price ranges (\$0-10, \$10-30, \$30-100, \$100+)
 - Different condition levels
 - Popular (many reviews) vs. Obscure (few reviews)

Example stratified evaluation:

```
python

    ○ # Group by book category
    ○ for category in ['Textbooks', 'Fiction', 'Non-Fiction', 'Collectibles']:
        subset = test_data[test_data['category'] == category]
        mae = mean_absolute_error(subset['actual'], subset['predicted'])
        r2 = r2_score(subset['actual'], subset['predicted'])

print(f"\n{category}: MAE=${mae:.2f}, R²={r2:.3f}, n={len(subset)}")
```

Typical findings:

- Textbooks: MAE = \$3.50, R² = 0.75 (predictable, condition-driven)
- Fiction: MAE = \$1.80, R² = 0.45 (lower prices, more variance in demand)
- Collectibles: MAE = \$45.00, R² = 0.30 (high variance, harder to predict)

Implication: May need separate models or confidence indicators per segment.

4.2 Training Procedures and Best Practices

Cross-Validation Strategy

Why cross-validation?

- Single train/test split can be lucky or unlucky
- Multiple folds provide robust performance estimates
- Essential for hyperparameter tuning (prevents overfitting to validation set)

K-Fold Cross-Validation (standard):

python

```
o from sklearn.model_selection import KFold
o kf = KFold(n_splits=5, shuffle=True, random_state=42)
o
o for fold, (train_idx, val_idx) in enumerate(kf.split(X)):
o     X_train, X_val = X[train_idx], X[val_idx]
o     y_train, y_val = y[train_idx], y[val_idx]
o
o     model.fit(X_train, y_train)
o     predictions = model.predict(X_val)
o
o     fold_rmse = sqrt(mean_squared_error(y_val, predictions))
o     print(f"Fold {fold+1} RMSE: ${fold_rmse:.2f}")
o
o # Average across folds
print(f"Mean CV RMSE: ${np.mean(fold_scores):.2f} ± ${np.std(fold_scores):.2f}")
```

Time-Based Split (for temporal data):

python

```
o # Simulate rolling forecast
o train_cutoff = '2024-06-01'
o val_cutoff = '2024-09-01'
o test_cutoff = '2024-12-01'
o
o train = data[data['date'] < train_cutoff]
o val = data[(data['date'] >= train_cutoff) & (data['date'] < val_cutoff)]
o
test = data[data['date'] >= val_cutoff]
```

Why time-based matters: Markets change. A model trained on 2020-2023 data predicting 2024 prices encounters:

- Different book releases
- Changed reader preferences (e.g., pandemic reading boom ending)
- Platform fee structure changes
- Economic conditions (inflation, recession)

Stratified sampling (for imbalanced categories):

```
python
    ○ from sklearn.model_selection import StratifiedKFold
    ○ skf = StratifiedKFold(n_splits=5)
    ○
    ○ # Stratify on price bins to ensure each fold has similar price distribution
    ○ price_bins = pd.qcut(y, q=5, labels=False) # Quintiles
    ○ for train_idx, val_idx in skf.split(X, price_bins):
# Train and validate
```

Hyperparameter Tuning

Grid Search (exhaustive):

```
python
    ○ from sklearn.model_selection import GridSearchCV
    ○
    ○ param_grid = {
    ○     'n_estimators': [100, 200, 500],
    ○     'max_depth': [5, 10, 15],
    ○     'learning_rate': [0.01, 0.05, 0.1],
    ○     'subsample': [0.8, 1.0]
    ○ }
    ○
    ○ grid_search = GridSearchCV(
    ○     XGBRegressor(),
    ○     param_grid,
    ○     cv=5,
    ○     scoring='neg_mean_squared_error',
    ○     n_jobs=-1
    ○ )
    ○
    ○ grid_search.fit(X_train, y_train)
```

```
best_model = grid_search.best_estimator_
```

Limitations: Computationally expensive (tries all combinations)

Random Search (more efficient):

```
python
o from sklearn.model_selection import RandomizedSearchCV
o
o param_distributions = {
o     'n_estimators': [100, 200, 300, 500, 1000],
o     'max_depth': range(3, 15),
o     'learning_rate': [0.001, 0.01, 0.05, 0.1, 0.2],
o     'subsample': [0.6, 0.7, 0.8, 0.9, 1.0],
o     'colsample_bytree': [0.6, 0.7, 0.8, 0.9, 1.0]
o }
o
o random_search = RandomizedSearchCV(
o     XGBRegressor(),
o     param_distributions,
o     n_iter=50, # Try 50 random combinations
o     cv=5,
o     scoring='neg_mean_squared_error',
o     n_jobs=-1,
o     random_state=42
)
)
```

Bayesian Optimization (smartest):

```
python
o from skopt import BayesSearchCV
o
o param_space = {
o     'n_estimators': (100, 1000),
o     'max_depth': (3, 15),
o     'learning_rate': (0.001, 0.3, 'log-uniform'),
o     'subsample': (0.5, 1.0)
o }
o
o bayes_search = BayesSearchCV(
o     XGBRegressor(),
o     param_space,
```

```
    ○ n_iter=50,  
    ○ cv=5,  
    ○ scoring='neg_mean_squared_error'  
)
```

How it works: Uses previous trial results to intelligently choose next hyperparameters (explores promising regions, avoids poor ones)

Winner: MachineHack competition participant achieved ~65% accuracy using **Bayesian optimization** for LightGBM tuning.

Preventing Overfitting

Signs of overfitting:

```
python  
    ○ train_rmse = 1.50 # Very low error on training  
    ○ val_rmse = 8.75 # Much higher on validation  
  
# Large gap indicates overfitting
```

Remedies:

1. Regularization (for boosting):

```
python  
    ○ model = XGBRegressor(  
    ○ reg_alpha=1.0, # L1 regularization  
    ○ reg_lambda=1.0, # L2 regularization  
    ○ max_depth=5, # Limit tree depth  
    ○ min_child_weight=3 # Require minimum samples per leaf  
)
```

2. Early Stopping:

```
python  
    ○ model.fit(  
    ○ X_train, y_train,  
    ○ eval_set=[(X_val, y_val)],  
    ○ early_stopping_rounds=50, # Stop if no improvement for 50 rounds  
    ○ verbose=False
```

)

3. Reduce Feature Dimensionality:

- Remove low-importance features
- Apply PCA or feature selection
- Use fewer text features (top 200 vs. all 5000 TF-IDF terms)

4. Increase Training Data:

- More samples reduce overfitting risk
- Augment data if possible (e.g., scrape more listings)

5. Simpler Model:

- Linear regression instead of deep neural network
- Fewer trees in Random Forest
- Shallower trees in boosting

6. Dropout / Data Augmentation (for neural networks):

python

- o model = Sequential([
- o Dense(128, activation='relu'),
- o Dropout(0.3), *# Randomly zero 30% of neurons during training*
- o Dense(64, activation='relu'),
- o Dropout(0.3),
- o Dense(1)

])

Handling Imbalanced Data

Problem: Most books are \$5-25, few are \$100+

- Model optimizes overall error, ignores rare expensive books
- Under-represents collectibles, textbooks in predictions

Solutions:

1. Stratified Sampling:

- Ensure train/val/test splits have similar price distributions

2. Class Weights (if framing as classification):

```

python

    ○ # Give higher weight to rare price ranges
    ○ class_weights = {0: 1.0, 1: 1.5, 2: 3.0} # More weight on expensive bins

model.fit(X, y, sample_weight=compute_sample_weight('balanced', y))

```

3. Oversampling Rare Segments:

```

python

    ○ # SMOTE for regression (synthetic samples)
    ○ from imblearn.over_sampling import SMOTE
    ○ sm = SMOTE()

X_resampled, y_resampled = sm.fit_resample(X, y_binned)

```

4. Separate Models per segment:

- Train one model for books < \$20
- Another for \$20-100
- Another for \$100+
- Route predictions based on book attributes

5. Weighted Loss Functions:

```

python

    ○ def weighted_mse(y_true, y_pred):
    ○     # Weight errors by price (higher price = more important to get right)
    ○     weights = np.log1p(y_true) / np.mean(np.log1p(y_true))
    ○     return np.mean(weights * (y_true - y_pred)**2)
    ○     ...
    ○
    ○     #### 4.3 Common Pitfalls and How to Avoid Them
    ○
    ○     ##### Pitfall 1: Feature Leakage (Revisited)
    ○
    ○     **Most insidious example**: Including "time-to-sale" as a feature
    ○     - Books that sold quickly might have been priced low (easier to sell)
    ○     - But you don't know time-to-sale until *after* the sale
    ○     - Model learns: "fast sale → low price", but can't use this at prediction time
    ○
    ○     **Detection**:
    ○     - If validation performance is suspiciously perfect ( $R^2 > 0.95$ ), investigate

```