# Forecasting Models in Machine Learning: Zero to Industry-Level Guide

## 1. What is Forecasting?

Forecasting means predicting future values using past data—especially important when data is collected in a time sequence (time series). Unlike just "predicting," forecasting focuses on trends, seasonality, and natural order in data.

## 2. Why Use Forecasting?

**Planning**: Staffing (hospitals, call centers), inventory (retail/wholesale), cash flow (finance).

**Optimization**: Reduce waste or shortage by knowing what's coming.

**Risk Reduction**: Make proactive decisions, anticipate risks.

**Automation & Data Product**: Real-time forecasting supports autonomous operations (smart grids, AI-powered trading).

## 3. When Is Forecasting Used?

Any situation requiring future estimates based on past patterns:

- **Retail**: Demand, sales, supply chain
- **Energy**: Consumption/load
- **Healthcare**: Admissions, resource allocation
- **Finance**: Stock/crypto prices, credit, insurance claims
- **Web/Apps**: Traffic load, user activity

## 4. Key Forecasting Concepts

- **Time Series**: Data points ordered by time (hourly, daily, monthly)
- **Target Variable**: What you want to forecast (e.g., daily sales)
- **Features/Exogenous Variables**: Context that influences the target (e.g., holidays, prices, temperature)
- **Lag**: Past values used as features
- **Trend**: Upward or downward tendency over time
- **Seasonality**: Regular repeating patterns (e.g., weekends, holidays)
- **Stationarity**: Properties don't change over time (often required for classic models)

# 5. Major Industry-Standard Forecasting Models

## A. Classic/Statistical Models

- **Naive**: Last observed value as forecast.
- **Moving Average**: Average of previous N steps.
- **ARMA/ARIMA**: Autoregressive (use past values), Moving Average, and Integrated (difference to remove trend).
- **SARIMA**: ARIMA with seasonality.
- **Exponential Smoothing/Holt-Winters**: Weight recent observations more; support trend and seasonality.

## B. Machine Learning Models

- **Linear Regression**: Use features (lag, calendar, exogenous) to predict future.
- **Tree Models (Random Forest, XGBoost, LightGBM)**: Can use complex dependencies, more than just time info.
- **Prophet (by Facebook)**: User-friendly—captures trend, seasonality, holidays.
- **Deep Learning**: LSTM, GRU, Transformer—handle long sequences, nonlinearities, many predictors.

---

# 6. Mathematical Formulation

## General TS Model

$$y_{t+h} = f(y_t, y_{t-1}, \ldots, y_{t-p}, X_t, X_{t-1}, \ldots) + \epsilon$$

- $y_{t+h}$: forecast at future time ($h$ steps ahead)
- $f$: the model (can be ARIMA, ML, neural net, etc)
- $X_t$: extra variables ("exogenous")
- $\epsilon$: error/noise

## Example: AR(1) Model

$$y_t = \phi y_{t-1} + \epsilon_t$$

## Prophet Model (Additive)

$$y(t) = g(t) + s(t) + h(t) + \epsilon_t$$

- $g(t)$: trend
- $s(t)$: seasonality
- $h(t)$: holiday effect

---

# 7. Complete Forecasting Pipeline

## Step 1: Data Understanding & Exploration

Load data, plot trends, detect missing, outliers.

**Example code:**
```
import matplotlib.pyplot as plt
plt.plot(df['date'], df['value'])
```

## Step 2: Feature Engineering

Create lag features, rolling means, categorize dates, include external predictors.

**Example code:**
```
df['lag_1'] = df['value'].shift(1)
df['rolling_7'] = df['value'].rolling(7).mean()
```

## Step 3: Data Splitting

Split by time order (train on past, test on future).

## Step 4: Model Selection & Fitting

**Classical: statsmodels (ARIMA, SARIMA)**
```
from statsmodels.tsa.arima.model import ARIMA
model = ARIMA(df['value'], order=(1,1,1))
result = model.fit()
```

**Prophet:**
```
from prophet import Prophet
df_prophet = df.rename(columns={'date':'ds','value':'y'})
m = Prophet()
m.fit(df_prophet)
```

**ML (Random Forest/XGBoost):**
```
from xgboost import XGBRegressor
model = XGBRegressor()
model.fit(X_train, y_train)
```

## Step 5: Prediction

Generate forecasts on test set (or future horizon)

**Example code:**
```
forecast = result.forecast(steps=10)
```

## Step 6: Evaluation (Key Metrics)

- **MAE**: Average absolute error
- **RMSE**: Penalizes large errors more
- **MAPE**: For percent errors (be careful with zeros)

**Example code:**
```
from sklearn.metrics import mean_absolute_error
mae = mean_absolute_error(y_test, y_pred)
```

Step 7: Visualization

Plot predictions vs actual

Plot residuals (should look random, not show trend)

---

# 8. When to Use Which Model?

- **Data is short, linear, or well-behaved:** ARIMA/SARIMA
- **Clear, strong seasonality or holidays:** Holt-Winters, Prophet
- **Multiple predictors, complex relationships, nonlinearities:** Tree models, XGBoost
- **Long sequences, lots of training data:** LSTM, GRU, Transformer

---

# 9. Final Best Practices

- Always hold out some "future" data for testing
- Feature engineering (lags, rolling stats, calendar, and exogenous variables) is often more important than the model used
- Visualize everything: trends, predictions, errors
- Watch out for data leakage (using future info to predict the past)

---

# Summary Table: Application Domains

| Domain | What is Forecasted | Typical Models |
|---|---|---|
| Retail | Sales/Demand | SARIMA, Prophet, XGBoost |
| Energy | Load, Consumption | ARIMA, LSTM, XGBoost |
| Finance | Stock/Price/Risk | ARIMA, LSTM, Random Forest |
| Healthcare | Admissions, Disease Spread | Prophet, SARIMA, ML |
| Web Traffic | Visits, Server Load | Prophet, LSTM, Tree Models |

Table 1: Application domains and typical forecasting models

---