# Digit Classification

## for Receipts and Invoices

Xingliang Shu

# Outline

1. Business Problem
2. Techniques to Classify Digits
   a. Dataset
   b. Data Visualization
   c. Models
   d. Comparison Table
   e. Misclassified data
3. Summary
4. Q&A
5. Appendix

# Business Problem
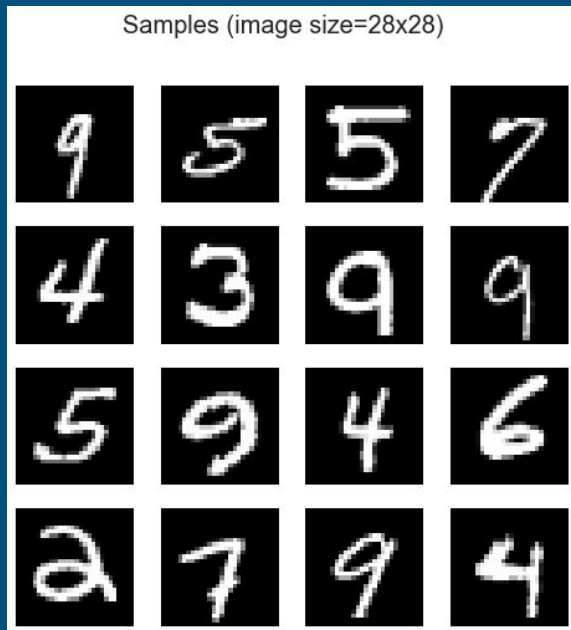
- Account Payable & Account Receivable
- Supply Chain Management

# Techniques to Classify Digits

# Dataset


Samples (image size=28x28)


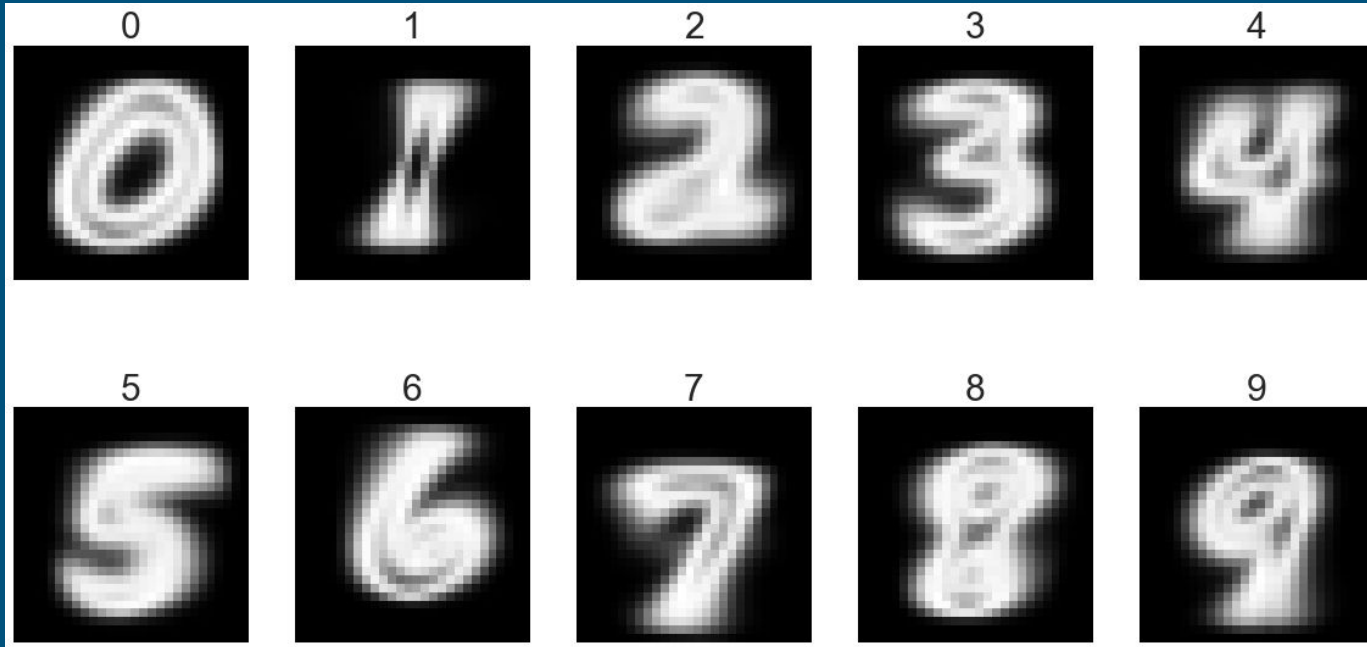Data Distribution

- 10 categories
- Training: 60k
- Testing: 10k
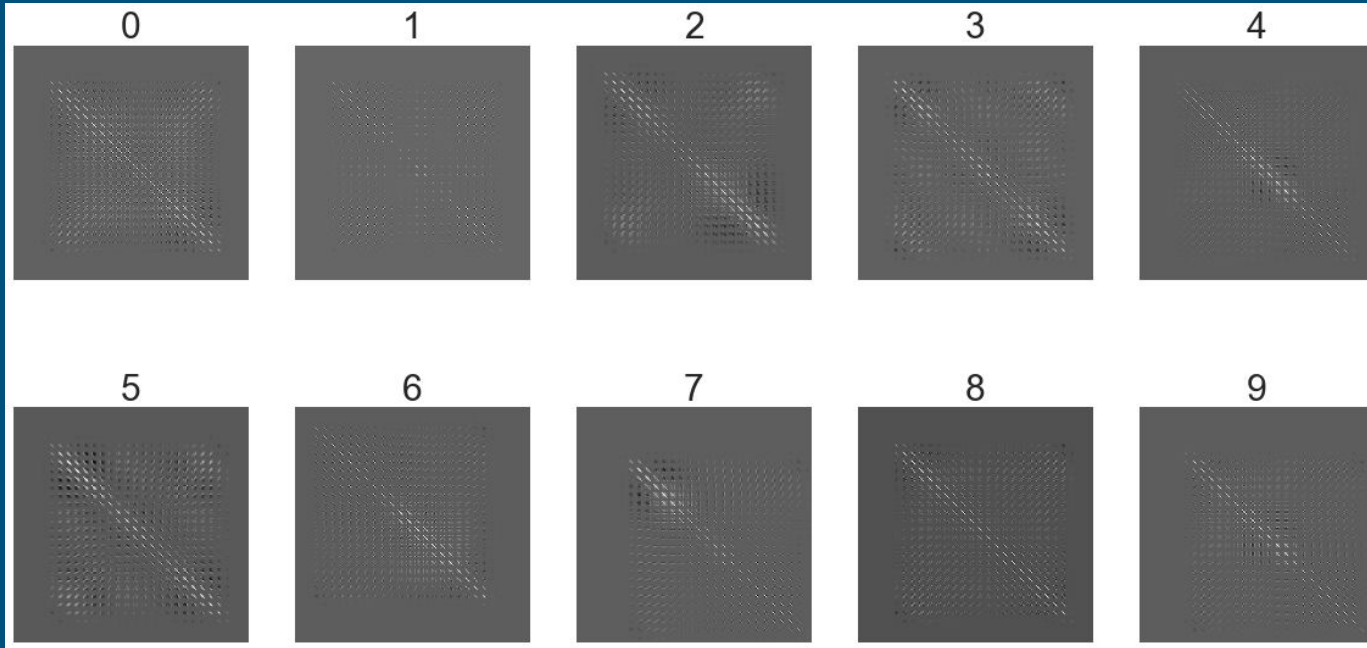- Metrice: accuracy

# Data Visualization (Means 2D)

# Data Visualization (Variances 2D)

# Data Visualization (Covariance Matrix)

# Multivariate Gaussian NB Version 1



Testing Dataset Result

- Parameters
  - Mean
  - Variance
- 4 ⇒ 9
- 5 ⇒ 3

# Multivariate Gaussian NB Version 2



- Parameters
  - Mean
  - Covariance matrix

# Binomial NB



Testing Dataset Result

- Parameters
  - Frequencies of 1
  - Frequencies of 0

# Comparison Table

| Model | Train acc | Val acc | Test acc |
|---|---|---|---|
| MVGNB 1 | 80.4% | 80.4% | 81.7% |
| BNB | 82.7% | 82.7% | 83.8% |
| **MVGNB 2** | **95.9%** | **95.3%** | **95.4%** |

- Best model: MVGNB 2

# Misclassified Data (2 ⇒ 7)


Classify 2 as 7

# Misclassified Data (3 ⇒ 5)



Classify 3 as 5

# Misclassified Data (5 ⇒ 3)



Classify 5 as 3

# Misclassified Data (4 ⇒ 9)



Classify 4 as 9

# Summary

- Usage of digit classifier
- Insights of parameters
- Best: MVGNB 2
- Misclassification
- Next
  - Digitize receipts
  - Study receipts
  - Advanced model
  - Automate digital process

# Q&A

# Appendix: for Technicians

# Data Visualization (Means 3D)

# Data Visualization (Variances 3D)

# Algorithm

- **What it does: given a sample predict its category.**
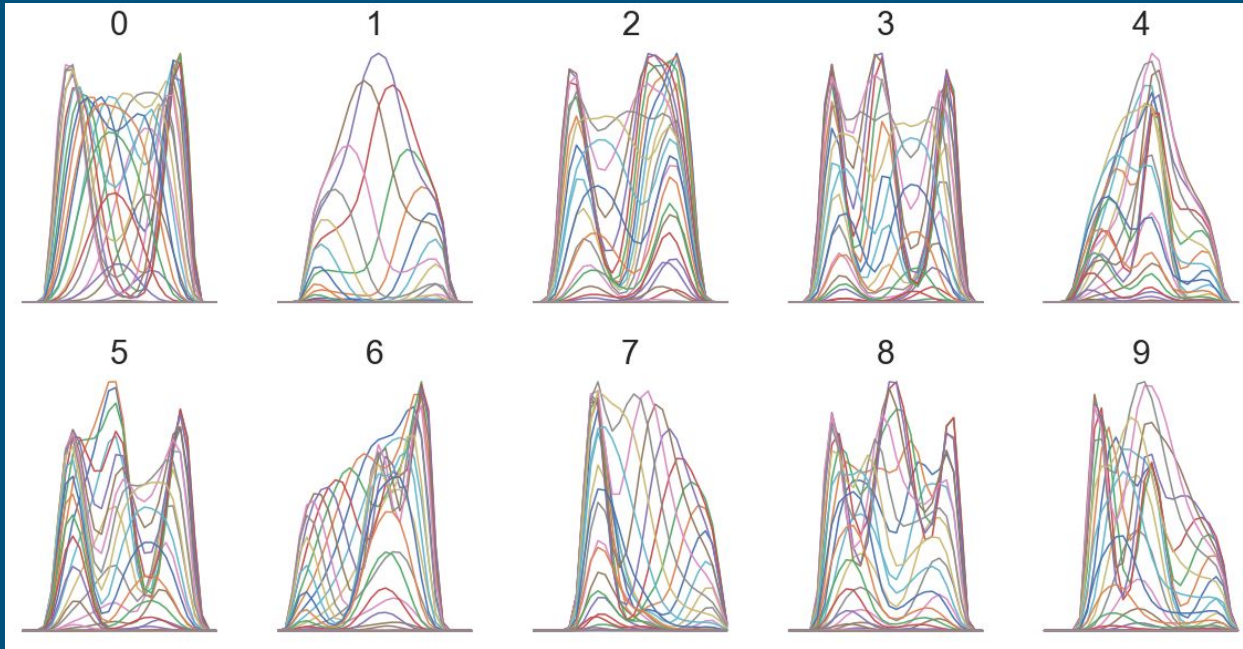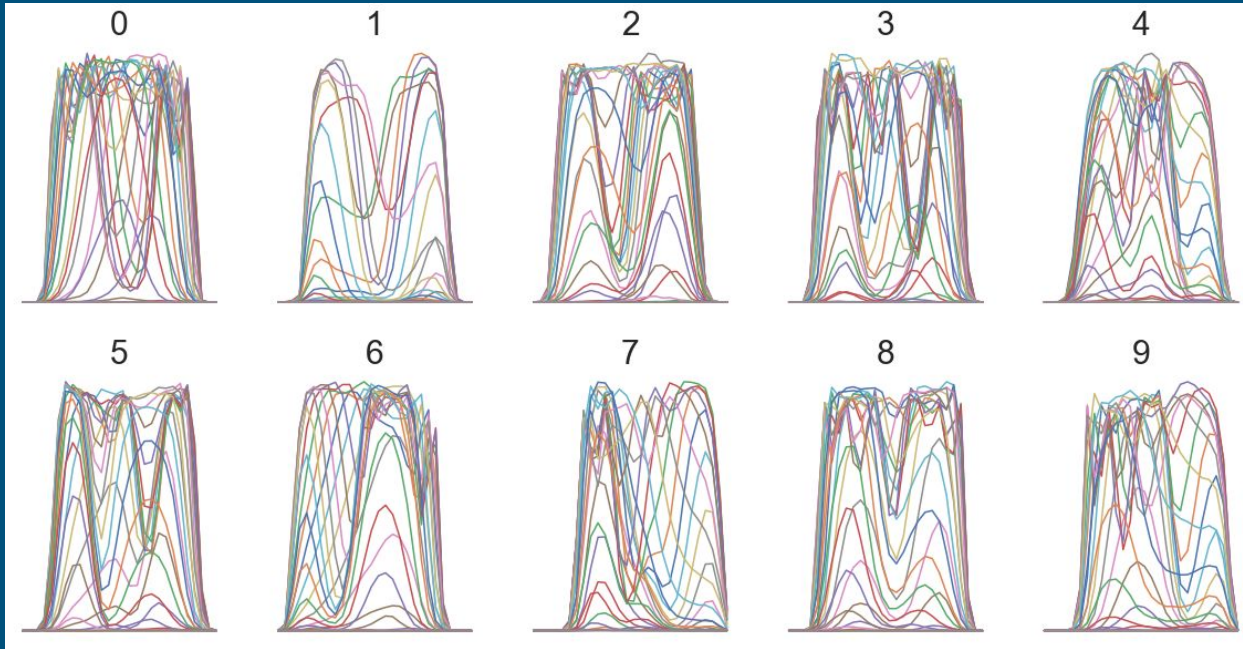- **$P(y|X) \propto argmax\{P(X|y)P(y)\}$** respect to y
- **X** represents a sample with **F** features, **F = [f1, f2, f3, f4, fF]**
- **y** represents a category.
- **P(y)** is the prior.
- Choose a model for **P(X|y)**: Gaussian, Binomial, etc
- For each category **k**:
  1. Calculate its **P(y)** where **P(y)=(#of samples in k)/(entire dataset)**
  2. Obtain **parameters** based on the selected model.
     a. If it's Gaussian, then ($\mu$ and $\sigma^2$) or ($\mu$ **and** $\Sigma$).
     b. If it's Binomial, then the **frequencies of ones** and **zeros**.
     c. Etc.

# Multivariate Gaussian Version 1

- How to predict?
  - Given a sample: X
  - X = [x1, x2, x3, …, xi, …]
  - $gi(xi) = (2*pi*\sigma^2)^{-1/2}exp(-0.5*(xi-\mu i)^2/\sigma^2)$
  - P_k = ∏gi * prior_k
  - P = [P_0, P_1, … P_k …, P_9]
  - argmax(P)

# Multivariate Gaussian Version 2

- How to predict?
    - Given a sample: X
    - X = [x1, x2, x3, …, xi, …]
    - For each category k
    - $g\_k(X) = \det(\Sigma\_k)^{-1/2}(2*pi*)^{-D/2}\exp(-0.5*(X-\mu\_k)^{T}\Sigma\_k(X-\mu\_k))$
    - P_k = g(X) * prior_k
    - P = [P_0, P_1, … P_k …, P_9]
    - argmax(P)

# Binomial

- Learned parameters: frequencies of 1s and 0s
- p denotes as the probability of 1
- $P(y|X) = p^X(1-p)^{(1-X)}$
- How to predict?
    - Given a sample: X
    - X = [0, 1, 1, 0, 1, … xi …, 1]
    - $P_k = p^X(1-p)^{(1-X)}*prior\_k$
    - P = [P_0, P_1, … P_k …, P_9]
    - argmax(P)