

CSE 5523: HW2



Outline

- You are to implement:
 - Pocket algorithm (improved perceptron)
 - Linear Gaussian discriminative analysis
 - Nonlinear Gaussian discriminative analysis

Data

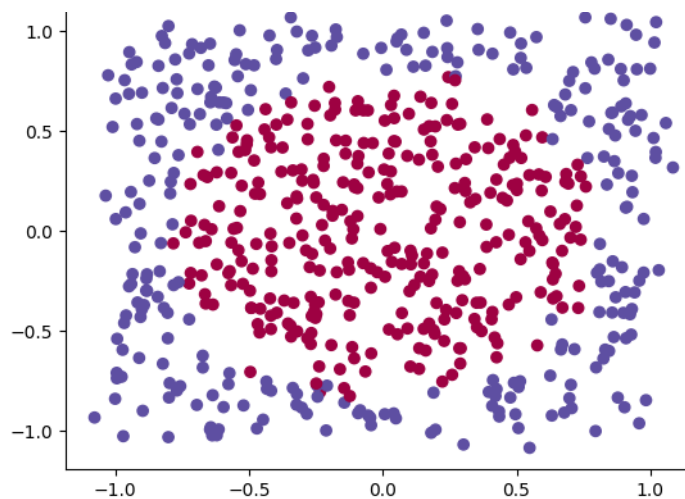
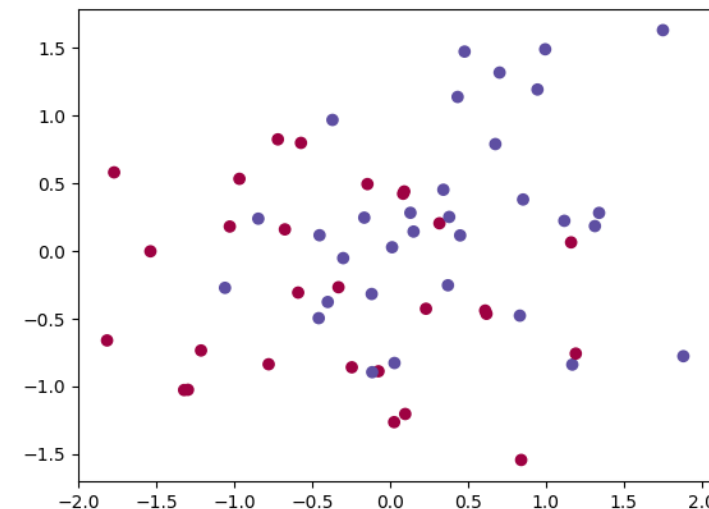
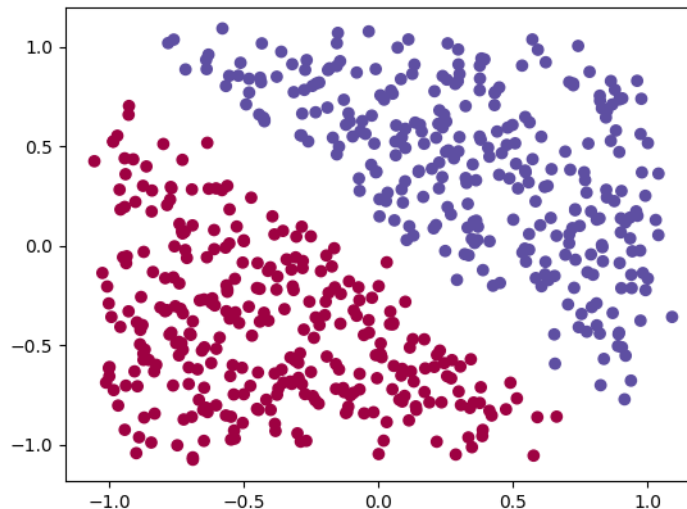
- Four data source

- 2D linear
- 2D noisy linear
- 2D quadratic (circle)
- MNIST (<5 vs. >=5)

- $X \in \mathbb{R}^{D \times N}$:

- A column as an instance

- $Y \in \{+1, -1\}^{N \times 1}$



Data

- The data \mathbf{X} are not appended with “1” yet.
- For feature transform for a 2D data instance $\mathbf{x} \in \mathbb{R}^2$, we do

$$\circ \boldsymbol{\phi}(\mathbf{x}) = \begin{bmatrix} x[1] \\ x[2] \\ x[1]^2 \\ x[2]^2 \\ x[1] \times x[2] \end{bmatrix}$$

- Again, you need to append “1” to the data $\boldsymbol{\phi}(\mathbf{x})$ if you want to solve $\tilde{\mathbf{w}}$ directly
- In the homework, we have done $\boldsymbol{\phi}(\mathbf{x})$ for you!

Accuracy

- Data = $\{ (\mathbf{x}_i, y_i \in \{+1, -1\}) \}_{i=1}^N$
- Accuracy = $\frac{1}{N} \sum_{i=1}^N \mathbf{1}[\hat{y}_i == y_i]$, where \hat{y}_i is the prediction based on \mathbf{x}_i

Pocket algorithm



Pocket algorithm

- **Training data:** $D_{tr} = \left\{ \left(\mathbf{x}_i \in \mathbb{R}^D, y_i \in \{+1, -1\} \right) \right\}_{i=1}^N$
- **Model:** $\text{sign}(\mathbf{w}\mathbf{x} + b) = \text{sign}(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}})$

Pocket algorithm

- Initialize $\tilde{\mathbf{w}}$ and $\tilde{\mathbf{w}}^{\text{best}} = \mathbf{0}$
- For $t = 1:T$
 - Loop for all training examples $\tilde{\mathbf{x}}_n$ (random order!)
 - Predict $\hat{y}_n = \text{sign}(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n)$
 - If $\hat{y}_n \neq y_n$
 - Update: $\tilde{\mathbf{w}} \leftarrow \tilde{\mathbf{w}} + \eta(y_n \tilde{\mathbf{x}}_n)$
 - Evaluate $\tilde{\mathbf{w}}$ on the “training data” and calculate the training accuracy
 - If training accuracy by $\tilde{\mathbf{w}}$ is “higher” than the training accuracy by $\tilde{\mathbf{w}}^{\text{best}}$
 - $\tilde{\mathbf{w}}^{\text{best}} \leftarrow \tilde{\mathbf{w}}$
- Output $\tilde{\mathbf{w}}^{\text{best}}$

See lecture 7 for extra details

Gaussian discriminant analysis



GDA

- **Training data:** $D_{tr} = \{ (\mathbf{x}_i \in \mathbb{R}^D, y_i \in \{+1, -1\}) \}_{i=1}^N$
- **Goal:** construct $p(Y = c|\mathbf{x})$ for $\hat{y} = \max_{c \in \{+1, -1\}} p(Y = c|\mathbf{x})$
- **Bayes' rules:** $p(Y = c|\mathbf{x}) \propto p(\mathbf{x}|Y = c)p(Y = c)$
 - $p(Y = c)$: Bernoulli
 - $p(\mathbf{x}|Y = c)$: multi-dimensional Gaussian

Nonlinear GDA

- $p(\mathbf{x} | Y = +1)$ and $p(\mathbf{x} | Y = -1)$ have their own covariance matrices Σ_{+1}, Σ_{-1}
- See slides 9, 10 for how to compute them
- See also your homework # 2

Linear GDA

- $p(\mathbf{x} | Y = +1)$ and $p(\mathbf{x} | Y = -1)$ share the same covariance matrix Σ
- Built upon the previous slide, given Σ_{+1}, Σ_{-1} and let N_{+1}, N_{-1} be the number of training examples per class, $\Sigma = \frac{N_{+1} \times \Sigma_{+1} + N_{-1} \times \Sigma_{-1}}{N}$
- See your homework # 2 for how to compute it

Prediction (please do “log” to prevent overflow)

$$\max_{c \in \{+1, -1\}} p(Y = c | \mathbf{x}) = \max_{c \in \{+1, -1\}} p(\mathbf{x} | Y = c) p(Y = c)$$

$$= \max_{c \in \{+1, -1\}} \log p(\mathbf{x} | Y = c) + \log p(Y = c)$$