

EECS 545: Machine Learning

Lecture 6. Classification 3

Honglak Lee
1/29/2025



Logistics/Announcements

- HW 2 due on Feb 11th
- Project Proposal due
 - No point deduction for HW1 through Jan 31th, Friday (3 late days with *consuming tokens* afterward)
 - 1 point per day from onwards

Last updated: 2/5/2025 12pm

2

Outline

(grey: already covered)

- Probabilistic discriminative models
 - ✓ Logistic Regression
 - ✓ Softmax Regression
- Probabilistic generative models
 - ✓ Gaussian discriminant analysis
 - Naive Bayes (continued)
- Discriminant functions (non-probabilistic)
 - Fisher's linear discriminant
 - Perceptron learning algorithm

3

Recap: Discriminative vs Generative Probabilistic Classifiers

- For classification, we want to compute $p(C_k | \mathbf{x})$
 - (a) **Discriminative** models: Directly model $p(C_k | \mathbf{x})$ and learn parameters from the training set.
 - Logistic regression
 - Softmax regression
 - (b) **Generative** models: Learn joint densities $p(\mathbf{x}, C_k)$ by learning $p(\mathbf{x} | C_k)$ and $p(C_k)$, and then use Bayes rule for predicting the class C_k given \mathbf{x} :
 - Gaussian Discriminant Analysis
 - Naive Bayes

4

Naive Bayes classifier (continued)

5

Naive Bayes classifier

- Probability of class label:
 - $p(C_k)$: Constant (e.g., Bernoulli)
- Conditional probability of data given the class
 - Naive Bayes assumption: $p(\mathbf{x} | C_k)$ is factorized (Each coordinate of \mathbf{x} is conditionally independent of other coordinates given the class label)

$$P(x_1, \dots, x_M | C_k) = P(x_1 | C_k) \cdots P(x_M | C_k) = \prod_{j=1}^M P(x_j | C_k)$$

- Classification: use Bayes rule

$$(\text{binary}) \quad P(C_1 | \mathbf{x}) = \frac{P(C_1, \mathbf{x})}{P(\mathbf{x})} = \frac{P(C_1, \mathbf{x})}{P(C_1, \mathbf{x}) + P(C_2, \mathbf{x})}$$

6

Naive Bayes classifier

- When classifying, we can simply find the class C_k that maximizes $P(C_k | \mathbf{x})$ using the Bayes rule:

$$\arg \max_k P(C_k | \mathbf{x}) = \arg \max_k P(C_k, \mathbf{x})$$

7

Naive Bayes classifier

- When classifying, we can simply find the class C_k that maximizes $P(C_k | \mathbf{x})$ using the Bayes rule:

$$\begin{aligned} \arg \max_k P(C_k | \mathbf{x}) &= \arg \max_k P(C_k, \mathbf{x}) \\ &= \arg \max_k P(C_k) P(\mathbf{x} | C_k) \end{aligned}$$

8

Naive Bayes classifier

- When classifying, we can simply find the class C_k that maximizes $P(C_k|\mathbf{x})$ using the Bayes rule:

$$\begin{aligned} \arg \max_k P(C_k|\mathbf{x}) &= \arg \max_k P(C_k, \mathbf{x}) \\ &= \arg \max_k P(C_k)P(\mathbf{x}|C_k) \\ \text{Naive Bayes assumption} &\quad \quad \quad = \arg \max_k P(C_k) \prod_{j=1}^M P(x_j|C_k) \end{aligned}$$

9

Example: Spam mail classification

- Label: $y=1$ (spam), $y=0$ (non-spam)
- Features:
 - x_j : j -th word in the mail, where M is the vocabulary size.
 - Multinomial variable (M -dimensional binary vector with only one coordinate with 1)
- Naive Bayes Assumption:
 - Given a class label y , each word in a mail is a independent multinomial variable.

10

Naive Bayes Spam classifier

- Model

$$\begin{aligned} P(\text{spam}) &= \text{Bernoulli}(\phi) \\ P(\text{word}|\text{spam}) &= \text{Multinomial}(\mu_1^s, \dots, \mu_M^s) \\ P(\text{word}|\text{nonspam}) &= \text{Multinomial}(\mu_1^{ns}, \dots, \mu_M^{ns}) \end{aligned}$$



11

Naive Bayes Spam classifier

- Model

$$\begin{aligned} P(\text{spam}) &= \text{Bernoulli}(\phi) \\ P(\text{word}|\text{spam}) &= \text{Multinomial}(\mu_1^s, \dots, \mu_M^s) \\ P(\text{word}|\text{nonspam}) &= \text{Multinomial}(\mu_1^{ns}, \dots, \mu_M^{ns}) \end{aligned}$$

- Goal

Find ϕ, μ^s, μ^{ns} that best fits the data $\{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$

12

Naive Bayes Spam classifier

- Model

$$\begin{aligned} P(\text{spam}) &= \text{Bernoulli}(\phi) \\ P(\text{word}|\text{spam}) &= \text{Multinomial}(\mu_1^s, \dots, \mu_M^s) \\ P(\text{word}|\text{nonspam}) &= \text{Multinomial}(\mu_1^{ns}, \dots, \mu_M^{ns}) \end{aligned}$$

- Goal

Find ϕ, μ^s, μ^{ns} that best fits the data $\{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$ by maximizing the joint likelihood:

$$\prod_{i=1}^N P(\mathbf{x}^{(i)}, y^{(i)})$$

- Joint Likelihood (joint probability of inputs/labels)
 - Note that the joint likelihood is conditioned on parameters ϕ, μ^s, μ^{ns}

13

Naive Bayes Spam classifier

- Model

$$\begin{aligned} P(\text{spam}) &= \text{Bernoulli}(\phi) \\ P(\text{word}|\text{spam}) &= \text{Multinomial}(\mu_1^s, \dots, \mu_M^s) \\ P(\text{word}|\text{nonspam}) &= \text{Multinomial}(\mu_1^{ns}, \dots, \mu_M^{ns}) \end{aligned}$$

- Goal

Find ϕ, μ^s, μ^{ns} that best fits the data $\{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$

- Likelihood - conditioned on parameters ϕ, μ^s, μ^{ns}

$$\begin{aligned} &\prod_{i=1}^N P(\mathbf{x}^{(i)}, y^{(i)}) \\ &= \prod_{i=1}^N P(\mathbf{x}^{(i)}|y^{(i)})P(y^{(i)}) \\ &= \left(\prod_{i:y^{(i)}=1} P(\mathbf{x}^{(i)}|y^{(i)})P(y^{(i)}) \right) \left(\prod_{i:y^{(i)}=0} P(\mathbf{x}^{(i)}|y^{(i)})P(y^{(i)}) \right) \end{aligned}$$

Spam Non-spam

14

Naive Bayes Spam classifier

- Likelihood - spam

$$\left(\prod_{i:y^{(i)}=1} P(\mathbf{x}^{(i)}|y^{(i)})P(y^{(i)}) \right) \quad \begin{array}{l} \nearrow \text{i-th mail} \\ \searrow \text{k-th word} \end{array}$$

- Naive Bayes assumption:

$$\begin{aligned} P(\text{spam}) &= \text{Bernoulli}(\phi) \\ P(\text{word}|\text{spam}) &= \text{Multinomial}(\mu_1^s, \dots, \mu_M^s) \\ \underline{P(\mathbf{x}^{(i)}|y^{(i)}=1)} &= \prod_{k=1}^{\text{len}(\mathbf{x}^{(i)})} P(x_k^{(i)}|y^{(i)}=1) \\ &= \prod_{k=1}^{\text{len}(\mathbf{x}^{(i)})} \prod_{j=1}^M (\mu_j^s)^{\mathbb{I}(x_k^{(i)} = \text{"j"th word})} \\ \underline{P(y^{(i)}=1)} &= \phi \end{aligned}$$

15

Naive Bayes Spam classifier

- Likelihood - spam (cont')

$$\begin{aligned} &\left(\prod_{i:y^{(i)}=1} P(\mathbf{x}^{(i)}|y^{(i)})P(y^{(i)}) \right) \\ &= \left(\prod_{i:y^{(i)}=1} \prod_{k=1}^{\text{len}(\mathbf{x}^{(i)})} \prod_{j=1}^M (\mu_j^s)^{\mathbb{I}(x_k^{(i)} = \text{"j"th word})} \phi \right) \end{aligned}$$

16

Naive Bayes Spam classifier

- Likelihood - spam (cont')

$$\begin{aligned} & \left(\prod_{i:y^{(i)}=1} P(\mathbf{x}^{(i)} | y^{(i)}) P(y^{(i)}) \right) \\ &= \left(\prod_{i:y^{(i)}=1} \prod_{k=1}^N \prod_{j=1}^M (\mu_j^s)^{\mathbb{I}(x_k^{(i)} = \text{"j" th word})} \phi \right) \\ &= \left(\prod_{i:y^{(i)}=1} \prod_{k=1}^N \prod_{j=1}^M (\mu_j^s)^{\mathbb{I}(x_k^{(i)} = \text{"j" th word})} \right) \left(\prod_{i:y^{(i)}=1} \phi \right) \end{aligned}$$

17

Naive Bayes Spam classifier

- Likelihood - spam (cont')

$$\begin{aligned} & \left(\prod_{i:y^{(i)}=1} P(\mathbf{x}^{(i)} | y^{(i)}) P(y^{(i)}) \right) \\ &= \left(\prod_{i:y^{(i)}=1} \prod_{k=1}^N \prod_{j=1}^M (\mu_j^s)^{\mathbb{I}(x_k^{(i)} = \text{"j" th word})} \phi \right) \\ &= \left(\prod_{i:y^{(i)}=1} \prod_{k=1}^N \prod_{j=1}^M (\mu_j^s)^{\mathbb{I}(x_k^{(i)} = \text{"j" th word})} \right) \left(\prod_{i:y^{(i)}=1} \phi \right) \\ &= \left(\prod_{j=1}^M (\mu_j^s)^{\sum_{i:y^{(i)}=1} \sum_{k=1}^N \mathbb{I}(x_k^{(i)} = \text{"j" th word})} \right) \left(\prod_{i:y^{(i)}=1} \phi \right) \end{aligned}$$

18

Naive Bayes Spam classifier

- Likelihood - spam (cont')

$$\begin{aligned} & \left(\prod_{i:y^{(i)}=1} P(\mathbf{x}^{(i)} | y^{(i)}) P(y^{(i)}) \right) \\ &= \left(\prod_{i:y^{(i)}=1} \prod_{k=1}^N \prod_{j=1}^M (\mu_j^s)^{\mathbb{I}(x_k^{(i)} = \text{"j" th word})} \phi \right) \\ &= \left(\prod_{i:y^{(i)}=1} \prod_{k=1}^N \prod_{j=1}^M (\mu_j^s)^{\mathbb{I}(x_k^{(i)} = \text{"j" th word})} \right) \left(\prod_{i:y^{(i)}=1} \phi \right) \\ &= \left(\prod_{j=1}^M (\mu_j^s)^{\sum_{i:y^{(i)}=1} \sum_{k=1}^N \mathbb{I}(x_k^{(i)} = \text{"j" th word})} \right) \left(\prod_{i:y^{(i)}=1} \phi \right) \\ &= \left(\prod_{j=1}^M (\mu_j^s)^{N_j^{spam}} \right) \phi^{N^{spam}} \end{aligned}$$

Definition:

N^{spam} : total # examples for spam

N_j^{spam} : total # of word j from the entire spam emails

19

Naive Bayes Spam classifier

- Likelihood - non-spam

$$\left(\prod_{i:y^{(i)}=0} P(\mathbf{x}^{(i)} | y^{(i)}) P(y^{(i)}) \right) \left(\prod_{i:y^{(i)}=0} \underline{P(\mathbf{x}^{(i)} | y^{(i)})} P(y^{(i)}) \right)$$

- Similarly for non-spam mails,

$$\begin{aligned} P(\mathbf{x}^{(i)} | y^{(i)} = 0) &= \prod_{k=1}^{\text{len}(\mathbf{x}^{(i)})} P(x_k^{(i)} | y^{(i)} = 0) \\ &= \prod_{k=1}^{\text{len}(\mathbf{x}^{(i)})} \prod_{j=1}^M (\mu_j^{ns})^{\mathbb{I}(x_k^{(i)} = \text{"j" th word})} \\ P(y^{(i)} = 0) &= 1 - \phi \end{aligned}$$

20

Maximum likelihood estimation

- Putting together:

$$\begin{aligned} & \prod_{i=1}^N P(\mathbf{x}^{(i)}, y^{(i)}) \\ &= \left(\prod_{i:y^{(i)}=1} P(\mathbf{x}^{(i)} | y^{(i)}) P(y^{(i)}) \right) \left(\prod_{i:y^{(i)}=0} P(\mathbf{x}^{(i)} | y^{(i)}) P(y^{(i)}) \right) \end{aligned}$$

21

Maximum likelihood estimation

- Putting together:

$$\begin{aligned} & \prod_{i=1}^N P(\mathbf{x}^{(i)}, y^{(i)}) \\ &= \left(\prod_{i:y^{(i)}=1} P(\mathbf{x}^{(i)} | y^{(i)}) P(y^{(i)}) \right) \left(\prod_{i:y^{(i)}=0} P(\mathbf{x}^{(i)} | y^{(i)}) P(y^{(i)}) \right) \\ &= \left(\phi^{N^{spam}} \prod_{word\ j} (\mu_j^s)^{N_j^{spam}} \right) \left((1 - \phi)^{N^{nonspam}} \prod_{word\ j} (\mu_j^{ns})^{N_j^{nonspam}} \right) \end{aligned}$$

Recall:

N^{spam} : total # examples for spam

$N^{nonspam}$: total # examples for non-spam

N_j^{spam} : total # word j from the entire spam emails

$N_j^{nonspam}$: total # word j from the entire nonspam emails

22

Maximum likelihood estimation

- Putting together:

$$\begin{aligned} & \prod_{i=1}^N P(\mathbf{x}^{(i)}, y^{(i)}) \\ &= \left(\prod_{i:y^{(i)}=1} P(\mathbf{x}^{(i)} | y^{(i)}) P(y^{(i)}) \right) \left(\prod_{i:y^{(i)}=0} P(\mathbf{x}^{(i)} | y^{(i)}) P(y^{(i)}) \right) \\ &= \left(\phi^{N^{spam}} \prod_{word\ j} (\mu_j^s)^{N_j^{spam}} \right) \left((1 - \phi)^{N^{nonspam}} \prod_{word\ j} (\mu_j^{ns})^{N_j^{nonspam}} \right) \end{aligned}$$

- Joint Log-likelihood

$$\begin{aligned} & \log P(\mathcal{D}) \\ &= \log \prod_{i=1}^N P(\mathbf{x}^{(i)}, y^{(i)}) \\ &= N^{spam} \log \phi + \sum_{word\ j} N_j^{spam} \log \mu_j^s + N^{nonspam} \log(1 - \phi) + \sum_{word\ j} N_j^{nonspam} \log \mu_j^{ns} \end{aligned}$$

23

Maximum likelihood estimation

- Joint Log-likelihood

$$\begin{aligned} & \log P(\mathcal{D}) \\ &= \log \prod_{i=1}^N P(\mathbf{x}^{(i)}, y^{(i)}) \\ &= N^{spam} \log \phi + \sum_{word\ j} N_j^{spam} \log \mu_j^s + N^{nonspam} \log(1 - \phi) + \sum_{word\ j} N_j^{nonspam} \log \mu_j^{ns} \end{aligned}$$

- Maximum-likelihood

– Take the derivative of log-likelihood w.r.t. the parameters, and set it to zero.

26

Maximum likelihood estimation

• From $\frac{\partial l}{\partial \phi} = \frac{1}{\phi} N^{spam} - \frac{1}{1-\phi} N^{nonspam} = 0$

We get $\phi = \frac{N^{spam}}{N^{spam} + N^{nonspam}}$

- Removing dependent variables:

$$\sum_{word=j=1}^M N_j^{spam} \log \mu_j^s = \sum_{word=j=1}^{M-1} N_j^{spam} \log \mu_j^s + N_M^{spam} \log(1 - \sum_{j=1}^{M-1} \mu_j^s)$$

$$\frac{\partial}{\partial \mu_j^s} \left(\sum_{word=j=1}^M N_j^{spam} \log \mu_j^s \right) = \frac{N_j^{spam}}{\mu_j^s} - \frac{N_M^{spam}}{1 - \sum_{j=1}^{M-1} \mu_j^s} = 0$$

s.t. $\sum_j \mu_j^s = 1$

$\sum_j \mu_j^{ns} = 1$

27

Maximum likelihood estimation

• From $\frac{\partial l}{\partial \phi} = \frac{1}{\phi} N^{spam} - \frac{1}{1-\phi} N^{nonspam} = 0$

We get $\phi = \frac{N^{spam}}{N^{spam} + N^{nonspam}}$

- Removing dependent variables:

$$\sum_{word=j=1}^M N_j^{spam} \log \mu_j^s = \sum_{word=j=1}^{M-1} N_j^{spam} \log \mu_j^s + N_M^{spam} \log(1 - \sum_{j=1}^{M-1} \mu_j^s)$$

$$\frac{\partial}{\partial \mu_j^s} \left(\sum_{word=j=1}^M N_j^{spam} \log \mu_j^s \right) = \frac{N_j^{spam}}{\mu_j^s} - \frac{N_M^{spam}}{1 - \sum_{j=1}^{M-1} \mu_j^s} = 0$$

$\frac{N_j^{spam}}{\mu_j^s} = \text{constant}, \forall j$

$\mu_j^s = \frac{N_j^{spam}}{\sum_j N_j^{spam}}$

28

Maximum likelihood estimation

- Summary:

$$P(spam) = \phi = \frac{N^{spam}}{N^{spam} + N^{nonspam}}$$

$$P(word = j | spam) = \mu_j^s = \frac{N_j^{spam}}{\sum_j N_j^{spam}}$$

$$P(word = j | non-spam) = \mu_j^{ns} = \frac{N_j^{nonspam}}{\sum_j N_j^{nonspam}}$$

Recall:

N^{spam} : total # examples for spam
 $N^{nonspam}$: total # examples for non-spam

N_j^{spam} : total # word j from the entire spam emails
 $N_j^{nonspam}$: total # word j from the entire nonspam emails

31

Laplace Smoothing

- Maximum likelihood is problematic when a specific word count is 0
 - Leads to probability of 0!
- Solution: Put “imaginary” counts for each word
 - prevent zero probability estimates (overfitting)!
 - E.g.: Adding “1” as imaginary count for each word

$$P(spam) = \phi = \frac{N^{spam}}{N^{spam} + N^{nonspam}}$$

$$P(word = j | spam) = \mu_j^s = \frac{N_j^{spam} + 1}{\sum_j N_j^{spam} + M}$$

$$P(word = j | non-spam) = \mu_j^{ns} = \frac{N_j^{nonspam} + 1}{\sum_j N_j^{nonspam} + M}$$

32

Outline

(grey: already covered)

- Probabilistic discriminative models
 - ✓ Logistic Regression
 - ✓ Softmax Regression
- Probabilistic generative models
 - ✓ Gaussian discriminant analysis
 - ✓ Naive Bayes
- Discriminant functions (non-probabilistic)
 - Fisher’s linear discriminant
 - Perceptron learning algorithm

33

Discriminant Functions

34

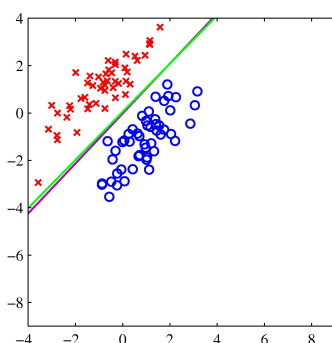
Linear Discriminant functions: Discriminating two classes

- Specify a weight vector \mathbf{w} and a bias w_0

$$h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

- Assign \mathbf{x} to C_1 if $h(\mathbf{x}) \geq 0$ and to C_0 otherwise.

- Q: How to pick \mathbf{w} ?



35

Linear Discriminant functions: Discriminating K>2 classes

- Instead each class C_k gets its own function

$$h_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k,0}$$

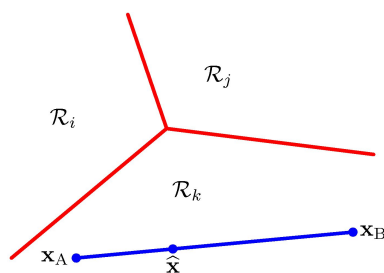
- Assign \mathbf{x} to C_k if

$$h_k(\mathbf{x}) > h_j(\mathbf{x}) \text{ for all } j \neq k$$

- The decision regions are convex polyhedra.

36

Decision Regions

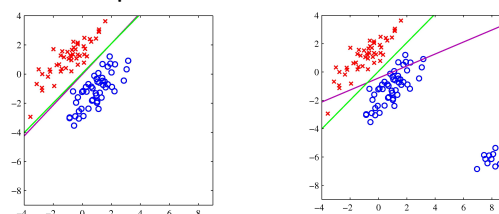


- Decision regions are convex, with piecewise linear boundaries.

37

How do we set the weights \mathbf{w} ?

- How about \mathbf{w} that minimizes squared error?
 - Label \mathbf{y} versus linear prediction $h(\mathbf{w})$.
 - Least squares is too sensitive to outliers. (why?)



Read Bishop book

38

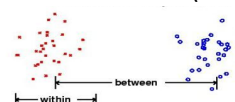
Learning Linear Discriminant Functions

- Fisher's linear discriminant
- Perceptron learning algorithm

39

Fisher's Linear Discriminant

- Let's consider binary classification case.
- Use \mathbf{w} to project \mathbf{x} to one dimension.
 - if $\mathbf{w}^T \mathbf{x} \geq -w_0$ then C_1 else C_0
- Select \mathbf{w} that best separates the classes.
- By "separating", the algorithm simultaneously
 - maximizes between-class (inter-class) variances
 - minimizes within-class (intra-class) variances

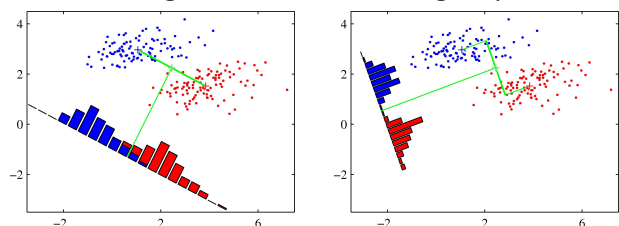


Read Bishop book

40

Fisher's Linear Discriminant

- Maximizing separation alone is not enough.
 - Minimizing class variance is a big help.



Read Bishop book

41

Objective function

- We want to maximize the "distance between classes"

$$\underline{m_2} - m_1 \equiv \mathbf{w}^T (\underline{\mathbf{m_2}} - \mathbf{m_1}) \quad \text{where } \mathbf{m}_k = \frac{1}{N_k} \sum_{n \in C_k} \mathbf{x}_n$$

Projected mean Mean

Read Bishop book

42

Objective function

- We want to maximize the "distance between classes"

$$\underline{m_2} - m_1 \equiv \mathbf{w}^T (\underline{\mathbf{m_2}} - \mathbf{m_1}) \quad \text{where } \mathbf{m}_k = \frac{1}{N_k} \sum_{n \in C_k} \mathbf{x}_n$$

Projected mean Mean

- While minimizing the "distance within each class"

$$s_1^2 + s_2^2 \equiv \sum_{n \in C_1} (\mathbf{w}^T \mathbf{x}_n - m_1)^2 + \sum_{n \in C_2} (\mathbf{w}^T \mathbf{x}_n - m_2)^2$$

Read Bishop book

43

Objective function

- We want to maximize the "distance between classes"

$$\underline{m_2} - m_1 \equiv \mathbf{w}^T (\underline{\mathbf{m_2}} - \mathbf{m_1}) \quad \text{where } \mathbf{m}_k = \frac{1}{N_k} \sum_{n \in C_k} \mathbf{x}_n$$

Projected mean Mean

- While minimizing the "distance within each class"

$$s_1^2 + s_2^2 \equiv \sum_{n \in C_1} (\mathbf{w}^T \mathbf{x}_n - m_1)^2 + \sum_{n \in C_2} (\mathbf{w}^T \mathbf{x}_n - m_2)^2$$

- Objective function: $J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}$

Read Bishop book

44

Derivation of objective

- Numerator: $m_2 - m_1 \equiv \mathbf{w}^\top (\mathbf{m}_2 - \mathbf{m}_1)$
 $\|m_2 - m_1\|^2 = \mathbf{w}^\top \underbrace{(\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^\top}_{= S_B} \mathbf{w}$

Read Bishop book
45

Derivation of objective

- Numerator: $m_2 - m_1 \equiv \mathbf{w}^\top (\mathbf{m}_2 - \mathbf{m}_1)$
 $\|m_2 - m_1\|^2 = \mathbf{w}^\top \underbrace{(\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^\top}_{= S_B} \mathbf{w}$
- Denominator:
 - $s_k^2 = \sum_{n \in C_k} (\mathbf{w}^\top \mathbf{x}_n - m_k)^2$
 $= \sum_{n \in C_k} \mathbf{w}^\top (\mathbf{x}_n - \mathbf{m}_k)(\mathbf{x}_n - \mathbf{m}_k)^\top \mathbf{w}$
 - $s_1^2 + s_2^2 = \mathbf{w}^\top \left[\sum_{k=1,2} \sum_{n \in C_k} (\mathbf{x}_n - \mathbf{m}_k)(\mathbf{x}_n - \mathbf{m}_k)^\top \right] \mathbf{w}$
 $= S_W$

Read Bishop book
46

Derivation of objective

- Numerator: $m_2 - m_1 \equiv \mathbf{w}^\top (\mathbf{m}_2 - \mathbf{m}_1)$
 $\|m_2 - m_1\|^2 = \mathbf{w}^\top \underbrace{(\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^\top}_{= S_B} \mathbf{w}$
- Denominator:
 - $s_k^2 = \sum_{n \in C_k} (\mathbf{w}^\top \mathbf{x}_n - m_k)^2$
 $= \sum_{n \in C_k} \mathbf{w}^\top (\mathbf{x}_n - \mathbf{m}_k)(\mathbf{x}_n - \mathbf{m}_k)^\top \mathbf{w}$
 - $s_1^2 + s_2^2 = \mathbf{w}^\top \left[\sum_{k=1,2} \sum_{n \in C_k} (\mathbf{x}_n - \mathbf{m}_k)(\mathbf{x}_n - \mathbf{m}_k)^\top \right] \mathbf{w}$
 $= S_W$
- After definition of terms, we get

$$J(\mathbf{w}) = \frac{\mathbf{w}^\top S_B \mathbf{w}}{\mathbf{w}^\top S_W \mathbf{w}}$$
 - Solution: $\mathbf{w} \propto S_W^{-1} (\mathbf{m}_2 - \mathbf{m}_1)$

Read Bishop book
47

Fisher's Linear Discriminant Analysis: Pros and Cons

Pros:

- Simple and effective approach for classification.
- Can effectively handle correlations between features
- Minimal assumptions about the underlying data distribution.
- Easy to interpret and explain

Cons:

- Only suitable for two-class classification problems
- Can be sensitive to outliers and may produce suboptimal results when the data has noisy features/labels

48

The Perceptron

- A "generalized linear function"

$$h(\mathbf{x}) = f(\mathbf{w}^\top \phi(\mathbf{x}))$$

where

$$f(a) = \begin{cases} +1, & a \geq 0 \\ -1, & a < 0 \end{cases}$$

- Uses target code: $y=+1$ for C_1 , $y=-1$ for C_2 .
- This means that we always want:

$$\mathbf{w}^\top \phi(\mathbf{x}^{(n)}) y^{(n)} > 0$$

49

The Perceptron Criterion

- Only count errors from misclassified points:

$$E_P(\mathbf{w}) = - \sum_{\mathbf{x}^{(n)} \in \mathcal{M}} \mathbf{w}^\top \phi(\mathbf{x}^{(n)}) y^{(n)}$$

– where \mathcal{M} is the set of **misclassified** points.

- Stochastic gradient descent:

– Update the weight vector according to the each misclassified sample (i.e., take gradient per sample):

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_P(\mathbf{w}) = \mathbf{w}^{(\tau)} + \eta \phi(\mathbf{x}^{(n)}) y^{(n)}$$

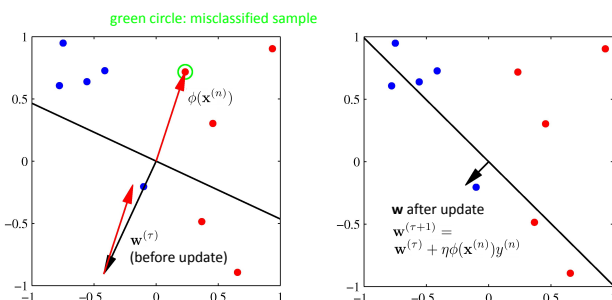
Note: update only for misclassified examples

50

Perceptron Learning (1)

- If $\mathbf{x}^{(n)}$ is misclassified, add $\phi(\mathbf{x}^{(n)})$ into \mathbf{w} .

red: positive ($y=+1$)
blue: negative ($y=-1$)

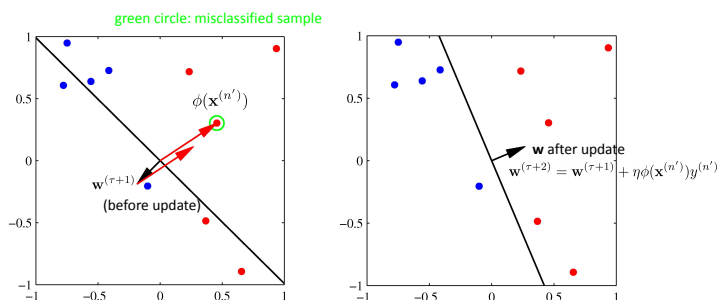


51

Perceptron Learning (2)

- If $\mathbf{x}^{(n)}$ is misclassified, add $\phi(\mathbf{x}^{(n)})$ into \mathbf{w} .

red: positive ($y=+1$)
blue: negative ($y=-1$)



52

Perceptron Learning

- Perceptron Convergence Theorem (Block, 1962, and Novikoff, 1962):
 - If there exists an exact solution (i.e., if the training data is linearly separable)
 - then the learning algorithm will find it in a finite number of steps.
- Limitations of perceptron learning:
 - The convergence can be very slow.
 - If dataset is not linearly separable, it won't converge.
 - Does not generalize well to $K > 2$ classes.

53

Next class

- Kernel methods
- Remember to submit your project proposals!

54

Any feedback (about lecture, slide, homework, project, etc.)?

(via anonymous google form: <https://forms.gle/fpYmi8tG9Me5qbP37>)



Change Log of lecture slides:

<https://docs.google.com/document/d/e/2PACX-1vSSIHIjklvpK7rkFSR1-5GYXyBCEW8UPtpSfCR9AR6M17K9ZOEmxfEwaWaW7kLDxusthsF8WICyZJ-/pub>

55