



CS 559: Machine Learning: Fundamentals & Applications

Week 5: Linear Classification II
Naïve Bayes Classifier & Logistic Regression
Spring 24





- 4.5. Probability Theory – MLE vs. MAP
- 4.6. Naïve Bayes Classifier
- 4.7. Logistic Regression
- 4.8. Model Selection - Regularization



4.5. Probability Theory – MLE vs. MAP

4.6. Naïve Bayes Classifier

4.7. Logistic Regression

4.8. Model Selection - Regularization

4.5. Probability Theory – MLE vs. MAP



- Maximum Likelihood Estimator (MLE)
 - Estimates the parameter, θ , of the possible probability that describes the data.
 - $x \sim p(x|\theta)$
- Maximum A Posteriori Estimation (MAP)
 - Treat the parameter, θ , as a random variable that describes the data.
 - $\underline{p(\theta|x)} \propto \underbrace{p(x|\theta)p(\theta)}_{\substack{\text{likelihood} \hookrightarrow \text{prior prob.} \\ \uparrow \text{posterior prob.}}}$ *assumptions*

$$\theta \uparrow \Rightarrow \text{Error} \downarrow$$

4.5. Probability Theory – MLE vs. MAP



- Example Coin Toss
 - Suppose a fair coin is tossed N times and each toss result is recorded.
 - Given data set:
 - $x = \{x_i\}_{i=1}^N, y_i \in \{-1, 1\}$
 - -1 is the tail and 1 is the head.
 - How can θ be estimated given x ?

$$\hat{\theta} = \frac{\sum_{i=1}^N I(y_i = 1)}{N}$$

↑ count heads *↑ total # of tosses*

4.5. Probability Theory – MLE vs. MAP



- MLE Approach
- Let's assume the probability of y is in the Bernoulli distribution.
- The Bernoulli distribution is the discrete probability distribution
 - $p(y=1) = p \quad \leftarrow \text{head}$
 - $p(y=0) = 1-p \quad \leftarrow \text{tail}$
 - $f(k|p) = \underline{p^k(1-p)^{1-k}}$

4.5. Probability Theory – MLE vs. MAP



- The example likelihood can be expressed as

$$p(x|\theta) = \theta^{n_1} (1-\theta)^{N-n_1}$$

$$y \sim p(x|\theta)$$

where n_1 is the count of $y=1$.

- Estimate $\hat{\theta}_{MLE}$: $\hat{\theta}_{MLE} = \frac{n_1}{N}$

$$\hat{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} \left(p(x|\theta) \right) = \underset{\theta}{\operatorname{argmax}} \left(\theta^{n_1} (1-\theta)^{N-n_1} \right)$$

$$= \underset{\theta}{\operatorname{argmax}} \left(\ln \left(\theta^{n_1} (1-\theta)^{N-n_1} \right) \right)$$

$$= \underset{\theta}{\operatorname{argmax}} \left[n_1 \ln \theta + (N-n_1) \ln (1-\theta) \right]$$

$$n_1 = \# \text{ of } y=1 \text{ (head)}$$

$$N = \text{total } \# \text{ of obs. (toss)}$$

$$n_1 = \frac{\sum I(y=1)}{N}$$

$$\theta^{n_1} = \text{total prob. of head}$$

$$(1-\theta)^{N-n_1} = \text{total prob. of tail.}$$

$$\ln(ab) = \ln a + \ln b$$

$$\ln(b^x) = x \ln b$$



$$\hat{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} \left(n_1 \ln \theta + (N - n_1) \ln (1 - \theta) \right)$$

take a derivative w.r.t. θ & set = 0.

$$= \underset{\theta}{\operatorname{argmax}} \left(\frac{n_1}{\theta} - \frac{N - n_1}{1 - \theta} \right)$$

= 0

$$\frac{n_1}{\theta} = \frac{N - n_1}{1 - \theta} \Rightarrow \hat{\theta}_{MLE} = \frac{n_1}{N}$$

$$f(\theta) = \ln \theta$$

$$f'(\theta) = \frac{df}{d\theta} = \frac{1}{\theta}$$

$$\frac{d}{d\theta} (\ln(1 - \theta))$$

$$= \frac{1}{1 - \theta} \frac{d}{d\theta} (1 - \theta)$$

$$= -\frac{1}{1 - \theta}$$

4.5. Probability Theory – MLE vs. MAP



- MAP approach
- Let θ is a random variable and estimate the posterior probability of θ given \mathbf{x} .
- Use the Bayes rule:

Posterior \propto likelihood * prior

$$p(\theta|\mathbf{x}) \propto \underbrace{p(\mathbf{x}|\theta)}_{\text{learning}} \underbrace{p(\theta)}_{\text{assumption}}$$

- The data set is $\mathbf{x} = \{y_i\}_{i=1}^N$.
- The actual Bayes rule is

$$p(\theta|\mathbf{x}) = \frac{p(\mathbf{x}|\theta)p(\theta)}{p(\mathbf{x})} \quad \begin{matrix} \leftarrow \text{actual distribution of data} \\ \text{very hard to know} \end{matrix}$$

$$p(\mathbf{x}) = \int p(\mathbf{x}|\theta)p(\theta) d\theta$$

4.5. Probability Theory – MLE vs. MAP



Steps to calculate $\hat{\theta}_{MAP}$:

- Prior probability: assume $p(\theta) \sim \theta^{\alpha-1} (1-\theta)^{\beta-1}$
- Likelihood: Bernoulli $p(x|\theta)$ [known] $N + n_1 \quad \theta^{n_1} (1-\theta)^{N-n_1} = p(x|\theta)$
- Posterior: compute $p(\theta)p(x|\theta)$
- Maximize the posterior.
- Solution of the example: $\hat{\theta}_{MAP} = \frac{(n_1 + \alpha - 1)}{N + \alpha + \beta + 2}$

$$\hat{\theta}_{MAP} = \arg \max_{\theta} (\theta^{n_1} (1-\theta)^{N-n_1} \cdot \theta^{\alpha-1} (1-\theta)^{\beta-1})$$

1. log-posterior

$$2. \frac{\partial}{\partial \theta} (\text{log-posterior}) = 0$$

3. Solve for θ .



$$P(\theta|x) = \left[\theta^{n_1} (1-\theta)^{N-n_1} \right] \left[\theta^{\alpha-1} (1-\theta)^{\beta-1} \right]$$

$$\begin{aligned} \ln(P(\theta|x)) &= n_1 \ln \theta + (N-n_1) \ln(1-\theta) + (\alpha-1) \ln \theta + (\beta-1) \ln(1-\theta) \\ &= (n_1 + \alpha - 1) \ln \theta + (N - n_1 + \beta - 1) \ln(1-\theta) \end{aligned}$$

$$\frac{\partial}{\partial \theta} \ln(P(\theta|x)) = \frac{n_1 + \alpha - 1}{\theta} - \frac{(N - n_1 + \beta - 1)}{1-\theta} = 0$$

$$\frac{\theta}{1-\theta} = \frac{n_1 + \alpha - 1}{N - n_1 + \beta - 1}$$

$$\text{if } N \rightarrow \infty \quad \lim_{N \rightarrow \infty} \left(\frac{n_1 + \alpha - 1}{N - n_1 + \beta - 1} \right) \approx \underbrace{\frac{n_1}{N - n_1}}_{\text{MLE}} \times \frac{\theta}{1-\theta} \Rightarrow \hat{\theta}_{\text{MAP}} = \frac{n_1}{N}$$

4.5. Probability Theory – MLE vs. MAP



- Example: Continuous $\mathbf{x} = \{x_i\}_{i=1}^N$, $x_i \in \mathbb{R}^d$. $\leftarrow \# \text{ of obs.}$ $\leftarrow \# \text{ of features}$
- Assume each x is independent. \rightarrow joint prob. of all features.
- The point distribution assumption:

$$x \sim \mathcal{N}(x|\mu, \sigma^2 = 1)$$

$$= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x - \mu)^2\right) \quad \begin{matrix} \text{Gaussian} \\ \hookdownarrow \end{matrix}$$

- Then the data distribution becomes

$$p(x|\mu) = \prod_{i=1}^N \frac{1}{(2\pi)^{\frac{d}{2}}} \exp\left(-\frac{1}{2}(x_i - \mu)^T(x_i - \mu)\right).$$

- The solution is $\hat{\mu}_{MLE} = \boxed{\frac{\sum_{i=1}^N x_i}{N}}$.

$$\frac{\partial}{\partial \mu} \ln p = 0$$

\rightarrow solve for μ ,

4.5. Probability Theory – MLE vs. MAP



- MLE vs. MAP:
- MAP works better when N is small, and the prior distribution is known.
- However, knowing the prior distribution for a data set with finite examples is difficult.
- MLE is preferred for large N .



4.5. Probability Theory – MLE vs. MAP

4.6. Naïve Bayes Classifier

4.7. Logistic Regression

4.8. Model Selection - Regularization

4.6. Naïve Bayes Classifier

posterior $\propto P(x|\omega)P(\omega)$



- The Naïve Bayes classifier is based on applying Bayes' theorem.
- Assumptions: Features are **conditionally independent** of a given target.
- Pros: Easy to train and predict fast. Sometimes, it is better than other linear classifiers.
- Cons: Strong assumptions.
 - It is not realistic in real life.
 - Computation becomes complex if training and test sets have different categories (We use Laplace estimation, see Bishop 4.4).

4.6. Naïve Bayes Classifier



- From Bayes' theorem, the conditional probability (posterior) can be decomposed as

$$p(C_k | \mathbf{x}) = \frac{p(\mathbf{x} | C_k)p(C_k)}{p(\mathbf{x})}$$

$$\frac{\prod_{k=1}^K p(C_k | \mathbf{x}_k)}{N} = p(C_k)$$

where C_k is the class and \mathbf{x} is the features.

- Remember that we are only interested in the numerator (=likelihood \times prior).
- The numerator is equivalent to the joint probability model using the chain rule.

$$p(y | \mathbf{x}) = \frac{p(x|y)p(y)}{p(x)}$$

4.6. Naïve Bayes Classifier



- When two variables, A and B, are independent,

$$p(A \perp B) = p(A)p(B)$$

$$\forall a, b: p(A = a \perp B = b) = p(A = a)p(B = b).$$

- If two variables, A and B, are **conditionally independent** given C,

$$p(A, B|C) = p(A|C)p(B|C)$$

$$\forall a, b, c: p(A = a \perp B = b | C = c) = p(A = a | C = c)p(B = b | C = c).$$

4.6. Naïve Bayes Classifier



- Suppose there are D-many features, $\mathbf{x} = [x_1, \dots, x_D]$.
- The joint probability is $p(C_k, x_1, \dots, x_n)$.
- Using the chain rule for the definition of conditional probability:

$$\begin{aligned} p(C_k, x_1, \dots, x_n) &= p(x_1, \dots, x_D, C_k) \\ &= p(x_1|x_2, \dots, C_k)p(x_2|x_3, \dots, C_k) \cdots p(x_D|C_k)p(C_k) \end{aligned}$$



4.6. Naïve Bayes Classifier

- Under the conditional independence assumption, the joint model can be expressed as

$$\begin{aligned} p(C_k | x_1, \dots, x_n) &\propto p(C_k, x_1, \dots, x_n) \\ &\propto p(C_k) \underbrace{p(x_1 | C_k) p(x_2 | C_k) \dots}_{n} \\ &\propto \underbrace{p(C_k)}_{\text{prior}} \prod_{i=1}^n \underbrace{p(x_i | C_k)}_{\text{likelihood}}. \end{aligned}$$

- The denominator $p(x)$ is

$$p(x) = \sum_k p(C_k) p(x | C_k).$$



4.6. Naïve Bayes Classifier

- Finding $p(C_k)$ is easy:

$$p(C_k = k) = \frac{\sum_{i=1}^N I(C_k = k)}{N} = \underline{\pi_k}$$

- The likelihood is

$$p(x|C_k) = \prod_{i=1} p(x_i|C_k)$$

$$p(C_k|x) \propto p(C_k) \prod_{i=1}^n p(x_i|C_k)$$

$$\ln(p(C_k|x)) = \ln p(C_k) + \ln \left(\prod_{i=1}^n p(x_i|C_k) \right)$$

4.6. Naïve Bayes Classifier



- The Bayes classifier is then

$$h(x) = \operatorname{argmax}_{C_k} \left(\sum_{i=1}^D \ln(p(x|C_k)) + \ln(p(C_k)) \right)$$

??

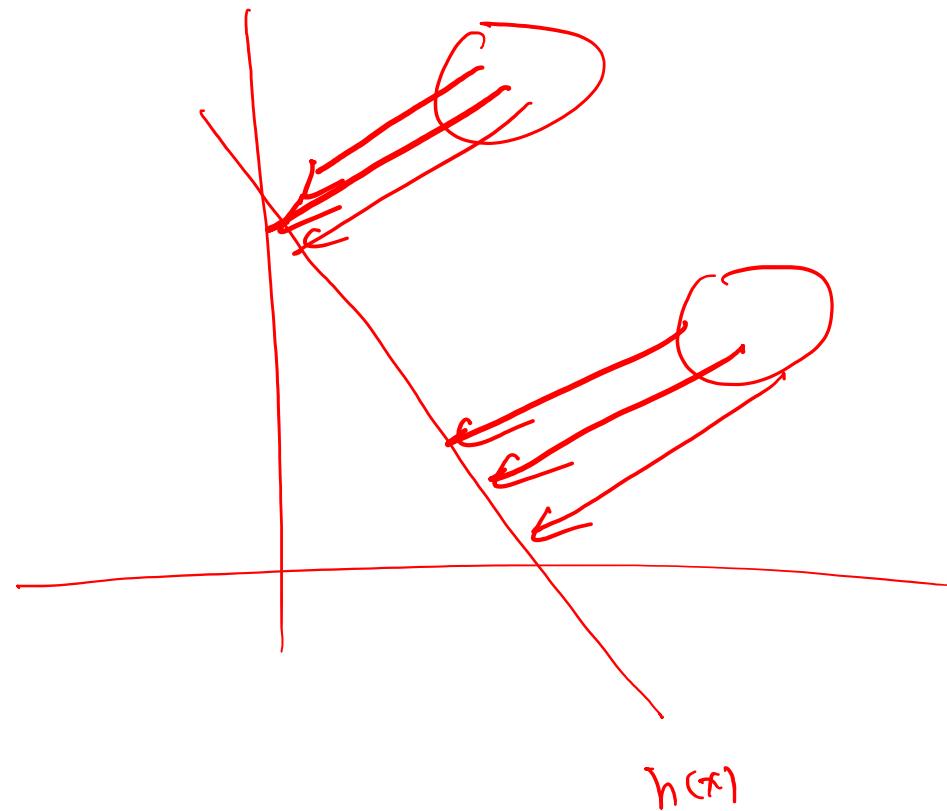
$$h(\vec{x}) = \omega^\top \vec{x} + \omega_0$$

4.6. Naïve Bayes Classifier



Assume features hold the LDA conditions. $\Sigma_1 = \Sigma_2 = \Sigma$

The Bayes classifier will show that the model is a linear equation.



$$\vec{\omega} \propto S^{-1} (\vec{m}_2 - \vec{m}_1)$$

$$\vec{\omega} \propto \Sigma^{-1} (\vec{\mu}_2 - \vec{\mu}_1)$$

4.6. Naïve Bayes Classifier



- Consider a binary class $y_i \in \{-1, 1\}$ and the continuous feature $\mathbf{x} \in \mathbb{R}^D$.
- Let the prior probability be $p(C_1) = \pi_1$ and $p(C_2) = \pi_2$. $= 1 - \pi_1$
- Assume \mathbf{x} is in Gaussian. $\mathbf{x} \sim \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$
- The probability distribution function (PDF) of the likelihood for class C_k is

$$p(\mathbf{x}|C_k) = \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\}. \quad \leftarrow$$

$$p(\vec{x}|C_1) = \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (\vec{x} - \vec{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\vec{x} - \vec{\mu}_1) \right)$$

4.6. Naïve Bayes Classifier



- The posterior becomes for C_1 :

$$P_{\text{Post}} = \frac{\text{Likelihood} \cdot \text{Prior}}{P(\vec{x})}$$

$$p(C_1|x) = \frac{p(x|C_1)p(C_1)}{p(x|C_1)p(C_1) + p(x|C_2)p(C_2)}$$

$$P(\vec{x}) = \sum_i P(C_i)P(x|C_i)$$

- Solutions:

$$\mathbf{w}^T = (\boldsymbol{\mu}_1^T - \boldsymbol{\mu}_2^T)\boldsymbol{\Sigma}^{-1}$$

$$w_0 = \frac{1}{2}(\boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2) + \ln \frac{\pi_1}{\pi_2}$$

?

$$\text{post.} = \frac{P(C_k) P(\vec{x}|C_k)}{P(\vec{x})}$$



$$P(C_1) = \pi_1 \quad P(C_2) = \pi_2$$

$$P(\vec{x}|C_1) = \text{const}_1 \cdot \exp\left(-\frac{1}{2}(\vec{x} - \vec{\mu}_1)^T \Sigma_1^{-1} (\vec{x} - \vec{\mu}_1)\right)$$

$$P(\vec{x}|C_2) = \text{const}_2 \cdot \exp\left(-\frac{1}{2}(\vec{x} - \vec{\mu}_2)^T \Sigma_2^{-1} (\vec{x} - \vec{\mu}_2)\right)$$

$$\Sigma_1 = \Sigma_2 = \Sigma$$

$\text{const}_1 = \text{const}_2 = \text{same}$

$$\frac{\pi_1 \cancel{\text{const}_1} \cdot \exp\left(-\frac{1}{2}(\vec{x} - \vec{\mu}_1)^T \underline{\Sigma_1^{-1}} (\vec{x} - \vec{\mu}_1)\right)}{\pi_1 \cancel{\text{const}_1} \cdot \exp\left(-\frac{1}{2}(\vec{x} - \vec{\mu}_1)^T \underline{\Sigma^{-1}} (\vec{x} - \vec{\mu}_1)\right) + \pi_2 \cancel{\text{const}_2} \cdot \exp\left(-\frac{1}{2}(\vec{x} - \vec{\mu}_2)^T \underline{\Sigma^{-1}} (\vec{x} - \vec{\mu}_2)\right)}$$

$$A = -\frac{1}{2}(\vec{x} - \vec{\mu}_1)^T \Sigma^{-1} (\vec{x} - \vec{\mu}_1)$$

$$B = -\frac{1}{2}(\vec{x} - \vec{\mu}_2)^T \Sigma^{-1} (\vec{x} - \vec{\mu}_2)$$



$$\text{post} = \frac{\pi_1 \exp(A)}{\pi_1 \exp(A) + \pi_2 \exp(B)} = \frac{1}{1 + \frac{\pi_2}{\pi_1} \frac{\exp B}{\exp A}}$$
$$= \frac{1}{1 + \frac{\pi_2}{\pi_1} \exp(B - A)}$$
$$= \frac{1}{1 + \exp(B - A + \ln \frac{\pi_2}{\pi_1})}$$

$$\ln(\exp(a)) = a$$



$$A = -\frac{1}{2}(\vec{x} - \mu_1)^T \Sigma^{-1} (\vec{x} - \mu_1) = -\frac{1}{2} \left[\underbrace{\vec{x}^T \Sigma^{-1} \vec{x}}_{\text{red wavy line}} - 2 \vec{x}^T \Sigma^{-1} \vec{\mu}_1 + \vec{\mu}_1^T \Sigma^{-1} \vec{\mu}_1 \right]$$

$$B = -\frac{1}{2} \left(\underbrace{\vec{x}^T \Sigma^{-1} \vec{x}}_{\text{red wavy line}} - 2 \vec{x}^T \Sigma^{-1} \vec{\mu}_2 + \vec{\mu}_2^T \Sigma^{-1} \vec{\mu}_2 \right)$$

$$B - A = \underbrace{\vec{x}^T \Sigma^{-1} \vec{\mu}_2}_{\text{blue underline}} - \frac{1}{2} \vec{\mu}_2^T \Sigma^{-1} \vec{\mu}_2 - \underbrace{\vec{x}^T \Sigma^{-1} \vec{\mu}_1}_{\text{blue underline}} + \frac{1}{2} \vec{\mu}_1^T \Sigma^{-1} \vec{\mu}_1$$

$$= - \left(\Sigma^{-1} \vec{x}^T (\vec{\mu}_1 - \vec{\mu}_2) \right) - \frac{1}{2} \left(\Sigma^{-1} (\vec{\mu}_2^T \vec{\mu}_2 - \vec{\mu}_1^T \vec{\mu}_1) \right)$$



Post =

$$\frac{1}{1 + \exp(B - A + \ln \pi_2 / \pi_1)}$$

$$B - A = -\Sigma^{-1}(\vec{\mu}_1 - \vec{\mu}_2)\vec{x}^T - \frac{1}{2}\Sigma^{-1}(\vec{\mu}_2^T \vec{\mu}_2 - \vec{\mu}_1^T \vec{\mu}_1)$$

post =

$$\frac{1}{1 + \exp\left(-\Sigma^{-1}(\vec{\mu}_1 - \vec{\mu}_2)\vec{x}^T - \frac{1}{2}\Sigma^{-1}(\vec{\mu}_2^T \vec{\mu}_2 - \vec{\mu}_1^T \vec{\mu}_1) + \ln \pi_2 / \pi_1\right)}$$

↑ const ↑ const ↑ Data ↑ const ↑ const.
 const const Data const const.
 ↗ w

$$= \frac{1}{1 + \exp\left(-\Sigma^{-1}(\vec{\mu}_1^T - \vec{\mu}_2^T)\vec{x} + \frac{1}{2}(\vec{\mu}_1^T \Sigma^{-1} \vec{\mu}_1 - \vec{\mu}_2^T \Sigma^{-1} \vec{\mu}_2) + \ln \pi_2 / \pi_1\right)}$$



$$\text{post} = \frac{1}{1 + \exp(-(\omega^T x + \omega_0))}$$

$$\omega^T = (\vec{\mu}_1^T - \vec{\mu}_2^T) \Sigma^{-1}$$

$$\omega_0 = \frac{1}{2} (\vec{\mu}_1^T \Sigma^{-1} \vec{\mu}_1 - \vec{\mu}_2^T \Sigma^{-1} \vec{\mu}_2) + \ln \frac{\pi_1}{\pi_2}$$

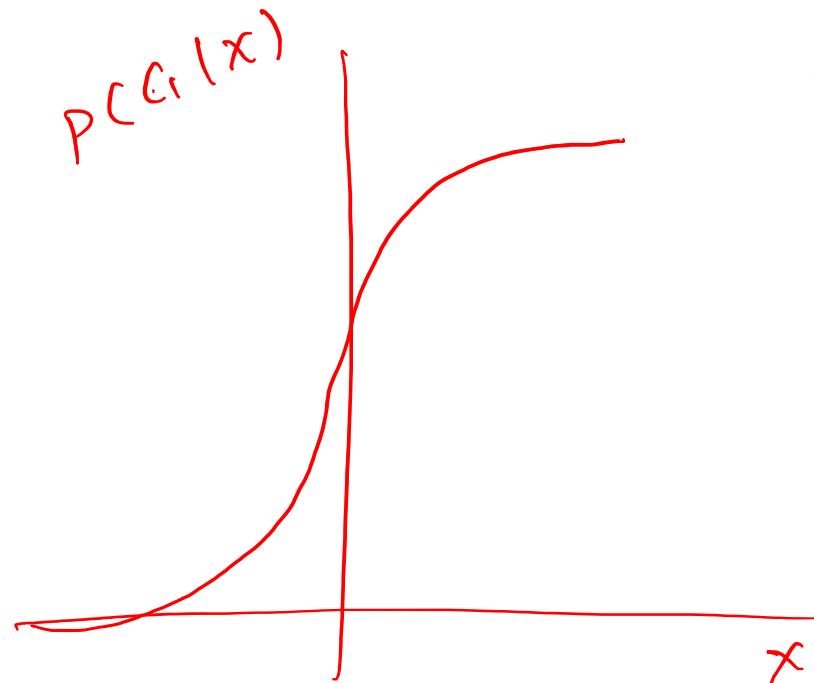
4.6. Naïve Bayes Classifier



- The solution, $\hat{h} = \underline{\mathbf{w}^T \mathbf{x} + w_0}$, and $p(C_1 | \mathbf{x})$ is called a **sigmoid function**.

$$p(C_1 | \mathbf{x}) = \frac{1}{1 + \exp(-\hat{h})}$$

- The predicted class for \mathbf{x} is



$$\hat{y} = \begin{cases} -1 & \text{if } \hat{h} < 0.5 \\ 1 & \text{otherwise} \end{cases}$$

LDA : $\hat{h} = \underline{\mathbf{w}^T \mathbf{x} + w_0}$

$$\hat{h} < 0 \quad \text{or} \quad \hat{h} \geq 0$$

$$C = -1 \quad C = 1$$

BC : $\frac{1}{1 + \exp(-\hat{h})} < 0.5 \quad \text{or}$

Prob. $\rightarrow \frac{1}{1 + \exp(-\hat{h})} \geq 0.5$

$$C = -1 \quad C = 1$$

4.6. Naïve Bayes Classifier



- Sigmoid function

$$f(x) = \sigma(x) = \frac{1}{1 + \exp(-x)}$$

4.6. Naïve Bayes Classifier



- How about MLE?
- Assume \mathbf{y} is in the Bernoulli distribution and \mathbf{X} is continuous data and is in Gaussian distribution.
- Each class has a Gaussian class-conditional density with a shared Σ .
- $y_n = \{1, 0\}$ for C_1 and C_2 , respectively.
- If the prior class probability is $p(C_1) = \pi$, $p(C_2) = 1 - \pi$.

4.6. Naïve Bayes Classifier



- For a data point x_i from class C_1 , $y_i = 1$ and hence the joint probability for each class becomes

$$p(x_i, C_1) = p(C_1)p(x_i|C_1) = \pi \mathcal{N}(x_i|\mu_1, \Sigma) \quad \leftarrow$$

$$p(x_i, C_2) = p(C_2)p(x_i|C_2) = (1 - \pi) \mathcal{N}(x_i|\mu_2, \Sigma). \quad \leftarrow$$

- Then the likelihood function is given by

$$p(X|y) = \prod_{i=1}^N [\pi \mathcal{N}(x_i|\mu_1, \Sigma)]^{y_i} [(1 - \pi) \mathcal{N}(x_i|\mu_2, \Sigma)]^{1-y_i}$$

$$\theta^{n_1} (1 - \theta)^{N - n_1}$$

4.6. Naïve Bayes Classifier



- Solutions:

$$\pi = \frac{1}{N} \sum_{i=1}^N y_i = \frac{N_1}{N}$$

$\frac{\partial \log(P(x|y))}{\partial \pi} \in \{1, 0\}$

where N_1 is the number of data points in class C_1 .

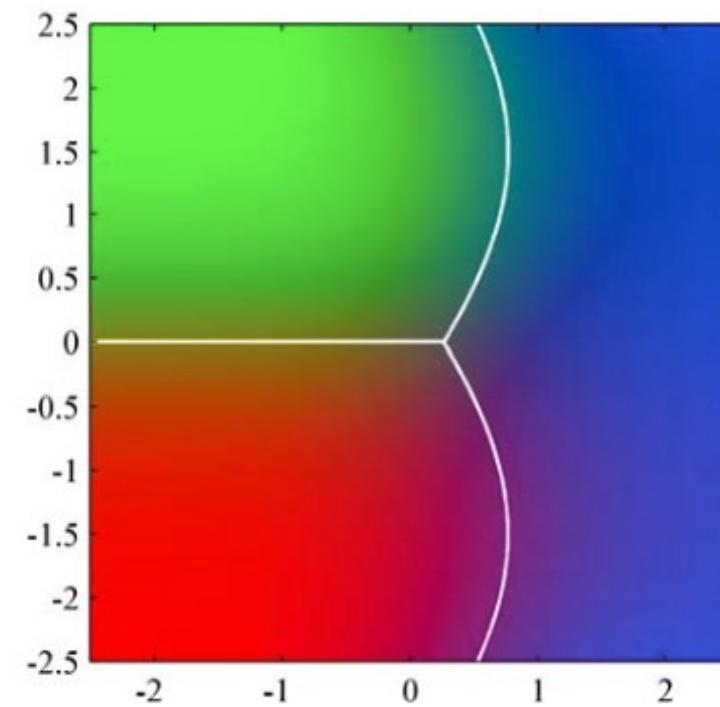
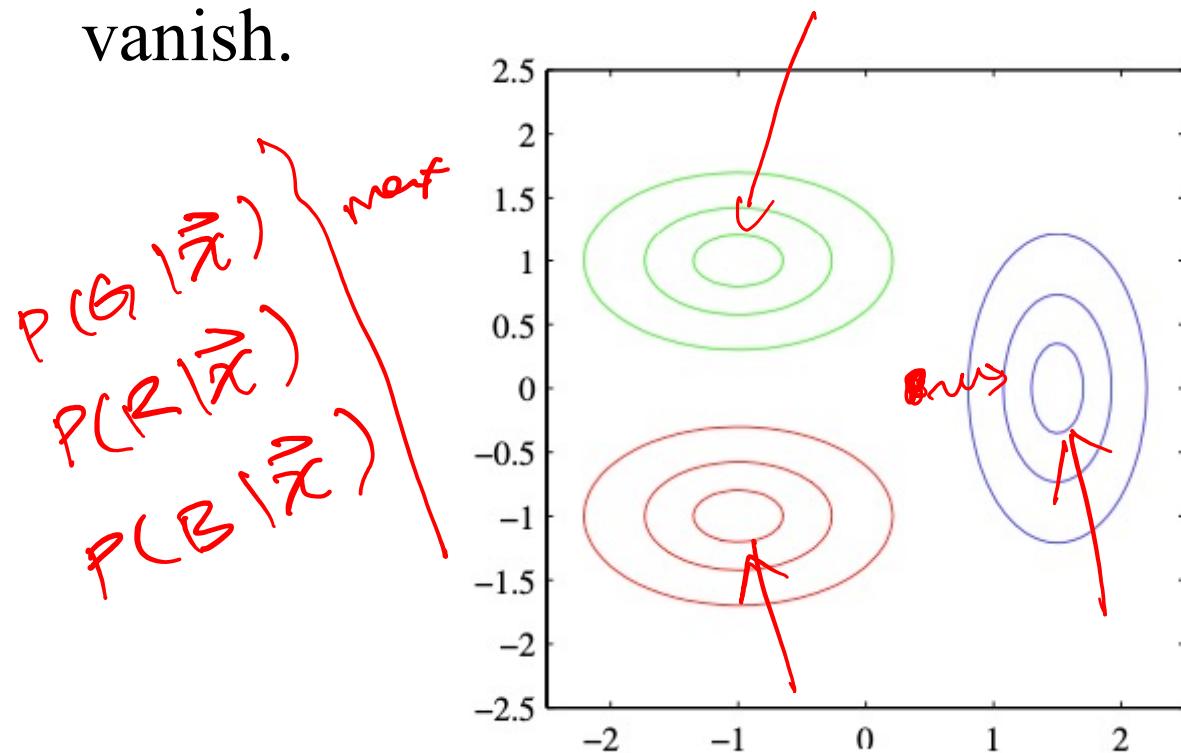
$$\mu_1 = \frac{1}{N_1} \sum_{i=1}^N y_i x_i, \mu_2 = \frac{1}{N_2} \sum_{i=1}^N (1 - y_i) x_i$$

$$\Sigma = \frac{N_1}{N} \mathbf{S}_1 + \frac{N_2}{N} \mathbf{S}_2$$

Recall, $w^T = (\mu_1^T - \mu_2^T) \Sigma^{-1}$.

4.6. Naïve Bayes Classifier

- For the case of $K > 2$ classes, $\hat{h}(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$
where $\underline{\mathbf{w}_k = \Sigma^{-1} \mu_k}$ and $\underline{w_{k0}} = -\frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \ln p(C_k)$.
- If each class has its own covariance matrix Σ_k , then $\underline{\mathbf{x} \Sigma_k^{-1} \mathbf{x}}$ term will not vanish.



4.6. Naïve Bayes Classifier



Class	(x_1, x_2)
1	(1,2), (2,3), (3,4.9)
0	(2,1), (3,2), (4, 3.9)

$$\Sigma = \begin{bmatrix} 0.667 & 0.967 \\ 0.967 & 1.467 \end{bmatrix}$$

$$\Sigma^{-1} = \begin{bmatrix} 48.22 & -32.22 \\ -32.22 & 22.22 \end{bmatrix}$$

$$w = (\mu_1^T - \mu_2^T) \Sigma^{-1} = [48.667, -80.444, 54.444]$$

$$\hat{w} = \frac{w}{|w|} = [0.448, -0.828, 0.560]$$

$$h = w^T x$$

$$y = \frac{1}{1 + \exp(-h)}$$

$$y \geq 0.5, \bar{y} = 1$$

$$y \leq 0.5, \bar{y} = 0$$

