



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.
 - $p(\theta|x) \propto \underbrace{p(x|\theta)p(\theta)}_{\substack{\text{likelihood} \\ \uparrow \text{posterior prob.}}} \text{ assumptions}$

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

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:
 - $\mathbf{x} = \{x_i\}_{i=1}^N, y_i \in \{-1, 1\}$
 - -1 is the tail and 1 is the head.
 - How can θ be estimated given \mathbf{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 $\hat{\theta}_{MLE}$: $\hat{\theta}_{MLE} = \frac{n_1}{N}$

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

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

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

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

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

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

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

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

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

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



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

take a derivative w.r.t. θ & set = 0.

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

$\underbrace{\qquad\qquad\qquad}_{=0}$

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

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

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

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

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

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

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

$$p(\theta|\mathbf{x}) \propto \underbrace{p(\mathbf{x}|\theta)}_{\text{not learning}} \underbrace{p(\theta)}_{\text{assumption}}$$

- 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 .



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

$$\ln(P(\theta|x)) = n_1 \ln \theta + (N-n_1) \ln(1-\theta) + (\alpha-1) \ln \theta + (\beta-1) \ln(1-\theta)$$

$$= (n_1 + \alpha - 1) \ln \theta + (N - n_1 + \beta - 1) \ln(1-\theta)$$

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

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

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



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:

$$x \sim \mathcal{N}(x|\mu, \sigma^2 = 1)$$
$$= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x - \mu)^2\right) \quad \leftarrow \text{Gaussian}$$

- 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

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

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

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.

$$p(y | \mathbf{x}) = \frac{p(x_1 | y)p(y)}{p(x)}$$

4.6. Naïve Bayes Classifier



$\perp\!\!\!\perp$ ← independent

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

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

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

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

$$p(A, B | C) = \boxed{p(A|C)} \boxed{p(B|C)}$$

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

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:

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

4.6. Naïve Bayes Classifier



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

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

joint P.
 posterior = $\frac{\text{likelihood} \times \text{prior}}{P(\vec{x})}$

- The denominator $p(x)$ is

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

joint probability of all classes

if a target is continuous

$$p(\vec{x}) = \int p(\vec{y}) p(\vec{x} | \vec{y}) d\vec{y}$$

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

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

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

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



4.6. Naïve Bayes Classifier

- The Bayes classifier is then

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

??

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

\uparrow

solution by opt. ω

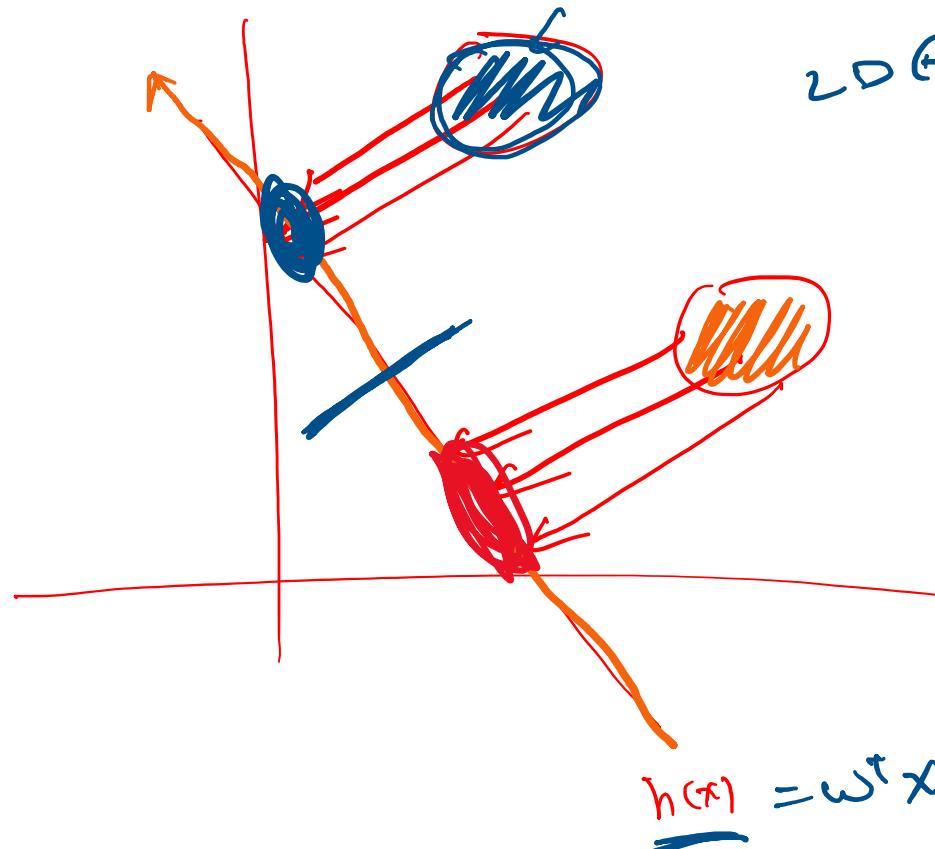
4.6. Naïve Bayes Classifier



Assume features hold the LDA conditions.

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

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



$$\text{LDA} \rightarrow \vec{\omega} \propto \underline{S^{-1}} (\vec{\mu}_2 - \vec{\mu}_1)$$
$$\vec{\omega} \propto \underline{\Sigma^{-1}} (\vec{\mu}_2 - \vec{\mu}_1) \in NBC.$$

↑
mean, sep.
inv. cov.

4.6. Naïve Bayes Classifier



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

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

on fire

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

$$p(\vec{x}|C_2) = \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (\vec{x} - \vec{\mu}_2)^T \Sigma^{-1} (\vec{x} - \vec{\mu}_2) \right)$$

4.6. Naïve Bayes Classifier



- The posterior becomes for C_1 :

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

likelihood prior
↓ ↘
 $P_{\text{POST}} = \frac{\text{Likelihood} \cdot \text{Prior}}{P(\vec{x})}$

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

- Solutions:

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

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

$$\text{post.} = \frac{P(C_k) P(\vec{x}|C_k)}{P(\vec{x})} \xrightarrow{\text{prior}} P(C_1) = \pi_1, P(C_2) = \pi_2$$

\vec{x}
likelihood

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

$$\pi_1 \cancel{\text{const}_1} \cdot \exp \left(-\frac{1}{2} (\vec{x} - \vec{\mu}_1)^T \cancel{\Sigma_1^{-1}} (\vec{x} - \vec{\mu}_1) \right) + \pi_2 \cancel{\text{const}_2} \cdot \exp \left(-\frac{1}{2} (\vec{x} - \vec{\mu}_2)^T \cancel{\Sigma_2^{-1}} (\vec{x} - \vec{\mu}_2) \right)$$

$$A = -\frac{1}{2} (\vec{x} - \vec{\mu}_1)^T \cancel{\Sigma_1^{-1}} (\vec{x} - \vec{\mu}_1) \leftarrow \text{class } 1$$

$$B = -\frac{1}{2} (\vec{x} - \vec{\mu}_2)^T \cancel{\Sigma_2^{-1}} (\vec{x} - \vec{\mu}_2) \leftarrow \text{class } 2$$



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

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

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

$$\frac{A}{A+B} = \frac{1}{1 + \frac{B}{A}}$$

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



$$A = -\frac{1}{2}(\vec{x} - \vec{\mu}_1)^T \Sigma^{-1} (\vec{x} - \vec{\mu}_1) = -\frac{1}{2} \left[\vec{x}^T \Sigma^{-1} \vec{x} \downarrow - 2 \vec{x}^T \Sigma^{-1} \vec{\mu}_1 + \vec{\mu}_1^T \Sigma^{-1} \vec{\mu}_1 \right]$$

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

$$\begin{aligned} B - A &= -\frac{1}{2} \left[\cancel{\vec{x}^T \Sigma^{-1} \vec{x}} - 2 \vec{x}^T \Sigma^{-1} \vec{\mu}_2 + \vec{\mu}_2^T \Sigma^{-1} \vec{\mu}_2 \right. \\ &\quad \left. - \cancel{\vec{x}^T \Sigma^{-1} \vec{x}} + \cancel{2 \vec{x}^T \Sigma^{-1} \vec{\mu}_1} - \vec{\mu}_1^T \Sigma^{-1} \vec{\mu}_1 \right] \\ &= -\frac{1}{2} \left[+2 \Sigma^{-1} (\vec{\mu}_1 - \vec{\mu}_2)^T \vec{x} + \Sigma^{-1} (\vec{\mu}_1^T \vec{\mu}_2 - \vec{\mu}_1^T \vec{\mu}_1) \right] \end{aligned}$$

$2 \vec{x}^T \Sigma^{-1} \vec{\mu}_1 - 2 \vec{x}^T \Sigma^{-1} \vec{\mu}_2$
 $2 [\vec{x}] \Sigma^{-1} \times (\vec{\mu}_1 - \vec{\mu}_2)$

Post =

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

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

post =

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

↑ const ↑ const Data const const.
↑ w w₀

$$= \frac{1}{1 + \exp(-(\underbrace{w^T x + w_0}_{\text{model}}))} = \frac{1}{1 + \exp(-h(x))}$$



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

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

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

4.6. Naïve Bayes Classifier



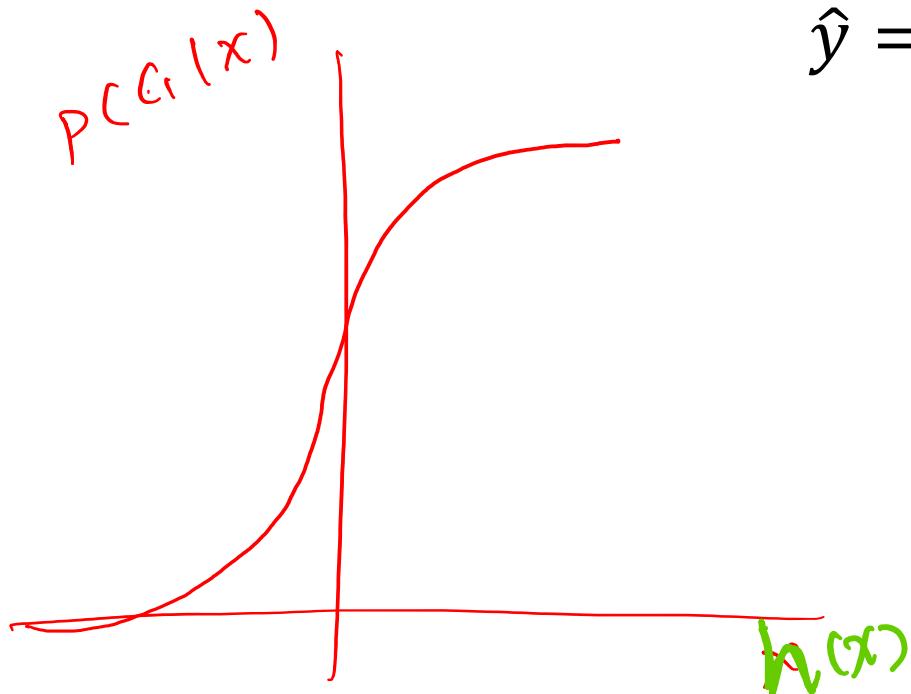
- The solution, $\hat{h} = \underline{\mathbf{w}^T \mathbf{x} + w_0}$, and $p(C_1 | \mathbf{x})$ is called a **sigmoid function**.

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

← posterior probability.

- The predicted class for \mathbf{x} is

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



$p(C_1 | \vec{x}) \in \{0, 1\}$.

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

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

$c = -1$ $c = 1$

BC

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

prob. \rightarrow

< 0.5
 $c = -1$

≥ 0.5
 $c = 1$

or \vec{x} has less than 50% of big class (-)



4.6. Naïve Bayes Classifier

- Sigmoid function

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

$$P(c_1 | \vec{x}) = \begin{cases} \text{posterior} \geq 0.5 & , c = -1 \\ < 0.5 & , c = 1 \end{cases}$$

$$P(c_1 | x_1) = 0.73$$

$$P(c_2 | x_1) = 1 - P(c_1 | x_1) = \underline{\underline{0.27}}$$

4.6. Naïve Bayes Classifier



- How about MLE? likelihood directly, no $p(\theta)$ assumption.
- 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

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

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

- Then the likelihood function is given by

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

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

$$\ln(p(x|y)) \rightarrow \frac{\partial}{\partial \mu_1} \ln(p(x|y)) = 0 \rightarrow \text{solve for } \underline{\mu_1}$$

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 .

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

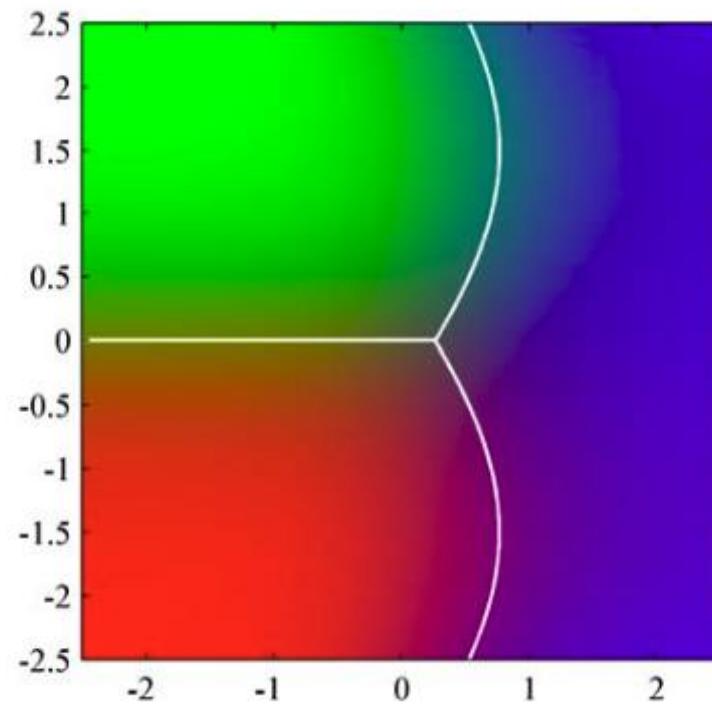
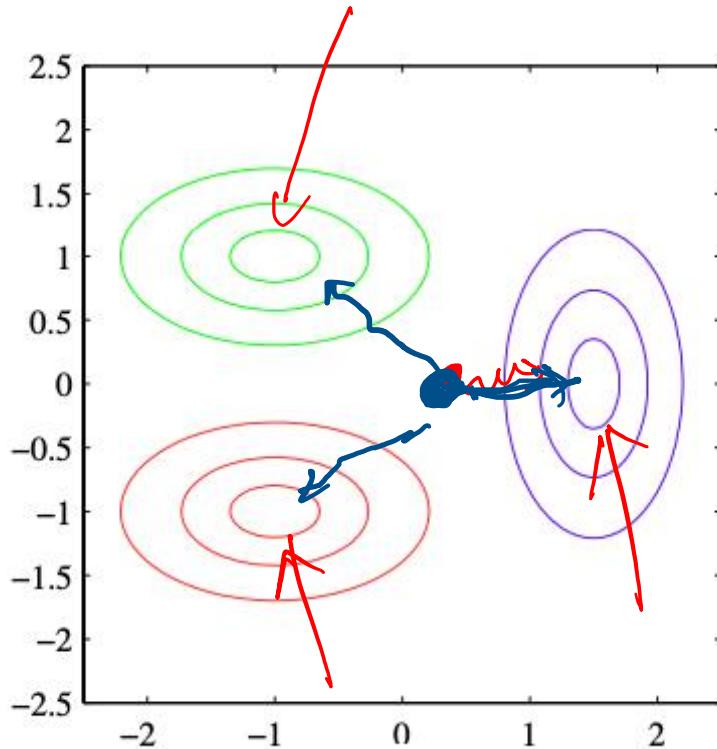
$$\Sigma = \frac{N_1}{N} \mathbf{S}_1 + \frac{N_2}{N} \mathbf{S}_2$$

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

$$\hat{y} = \max(P(G|\vec{x}), P(R|\vec{x}), P(B|\vec{x}))$$



4.6. Naïve Bayes Classifier



$$\hat{\mu}_1 = \frac{\sum_{i=1}^N x_i (class=1)}{N} = \frac{I(y=C_k)}{N}$$

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

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

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

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

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

$$\sigma(x) = \frac{1}{1 + \exp(0.560(x_2) + 0.828(x_1) + 0.448)}$$

