

Lecture 16.1

$$h(\vec{x}) = \arg \max_y P(\vec{x} | y) P_\theta(y)$$

$$= \prod_{d=1}^D P_\theta(x_d | y) = \frac{m!}{x_1! x_2! \dots x_d!} \prod_{d=1}^D ([\theta_y]_d)^{x_d}$$

$y \in \{\text{spam} | \text{ham}\}$

$$\frac{P_\theta(x | y=\text{spam}) P_\theta(y=\text{spam})}{P_\theta(x | y=\text{ham}) P_\theta(y=\text{ham})} = \frac{\cancel{m!} \left(\prod_{d=1}^D ([\theta_{\text{spam}}]_d)^{x_d} \right) P_\theta(y=\text{spam})}{\cancel{m!} \left(\prod_{d=1}^D ([\theta_{\text{ham}}]_d)^{x_d} \right) P_\theta(y=\text{ham})}$$

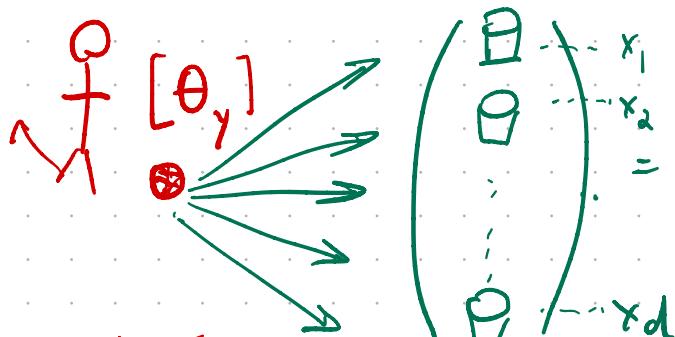
If the denominator \geq the divisor, output "spam".
O.t.w., output "ham".

Alternatively $\ln P_\theta(y|x) \propto \sum_{d=1}^D x_d \log([\theta_y]_d) + \log(P_\theta(y))$

Remark: The multinomial distribution is our modeling assumption for $P(x|y)$

$$P_\theta(x|y) = \frac{m!}{x_1! x_2! \dots x_d!} \prod_{d=1}^D ([\theta_y]_d)^{x_d}$$

multinomial distribution



Multinomial Feature

x_d - # of times the ball goes to bin d

$\in \{0, 1, \dots, m\}$

and $\sum_{d=1}^D x_d = m$

$$x = \begin{pmatrix} 2 \\ 1 \\ 2 \end{pmatrix}$$

x_1, x_2, x_3

x_1, x_2, x_3, x_4, x_5

Continuous Features (Gaussian Naive Bayes)

$x_\alpha \in \mathbb{R}$ (e.g. heights, weights)

recall our task to classify fish type
given the length

- Our modeling assumption:

- $\underset{\theta}{P}(\vec{x}|y) = \text{multivariate Gaussian}$

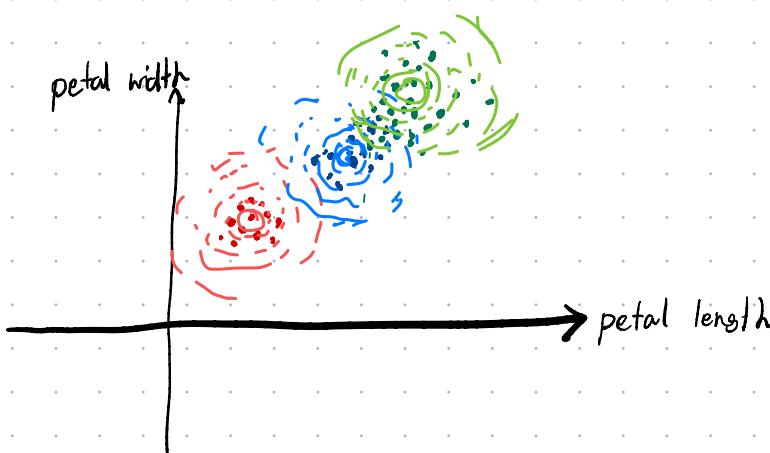
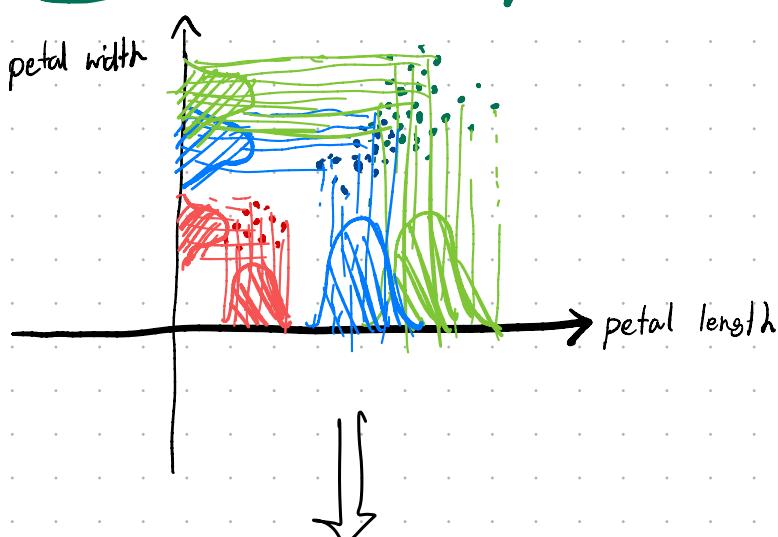
- $P_{\theta_\alpha}(\vec{x}_\alpha|y) = \text{univariate Gaussian}$
with some mean and s.d.

$$\stackrel{d}{=} N([\mu_y]_\alpha, [\sigma_y]_\alpha)$$

$\Rightarrow \prod_{\alpha=1}^d P_{\theta_\alpha}(\vec{x}_\alpha|y)$ is a multivariate Gaussian

[Gaussian \times Gaussian \rightarrow Gaussian]

Remark: We already know how to estimate $P(\vec{x}_\alpha|y)$



Summary of Naïve Bayes

- Naïve Bayes = Bayes classifier + naïve Bayes assumption
- The assumption says "all feature values are independent."
- We may have data that violates the assumptions.
- If your data is multinomial Features for binary classification, then naïve Bayes always give the linear decision boundary (naïve Bayes is a linear classifier)

- Assume $y \in \{-1, +1\}$

$$h(\vec{x}) = +1 \text{ iff } P(y=+1 | \vec{x}) > P(y=-1 | \vec{x})$$

$$\text{iff } P(\vec{x} | y=+1) P(y=+1) > P(\vec{x} | y=-1) P(y=-1)$$

$$\text{iff } \prod_{d=1}^D P(x_d | y=+1) P(y=+1) > \prod_{d=1}^D P(x_d | y=-1) P(y=-1)$$

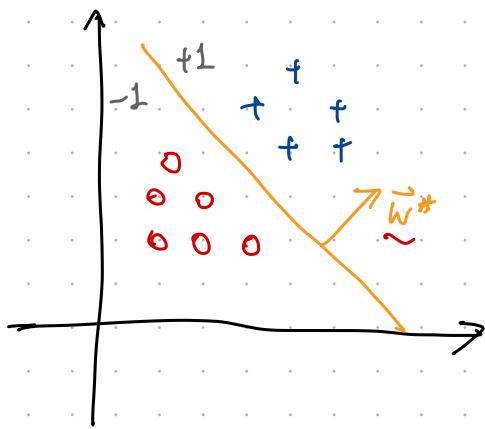
$$\text{iff } \sum_{d=1}^D x_d \log(P(x_d | y=+1)) + \log(P(y=+1)) > \sum_{d=1}^D x_d \log(P(x_d | y=-1)) + \log(P(y=-1))$$

$$\text{iff } \sum_{d=1}^D x_d \left(\log(P(x_d | y=+1)) - \log(P(x_d | y=-1)) \right) + \left(\log(P(y=+1)) - \log(P(y=-1)) \right) > 0$$

$$\text{iff } \sum_{d=1}^D x_d \underbrace{\left(\log([\theta_{y=+1}]_d) - \log([\theta_{y=-1}]_d) \right)}_{W} + \underbrace{\left(\log(P(y=+1)) - \log(P(y=-1)) \right)}_{b} > 0$$

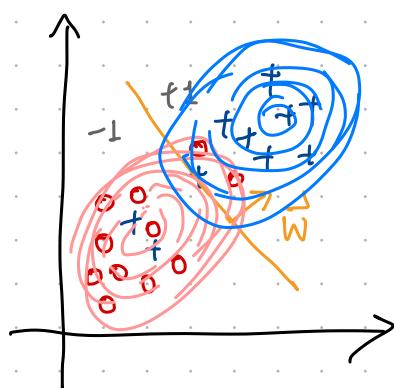
$$\text{iff } w^T x + b > 0$$

hyperplane



Perceptron

find \vec{w}^* that best separates negative points from positive points



Naïve Bayes + multinomial features

find \vec{w} that best separates the two modeling distributions

Discriminative algorithm: model $P(y|x)$ (e.g. k-NN, Perceptron)

Generative algorithm: model $P(x|y)$ and $P(y)$ to estimate $P(y|x)$ (e.g. naïve Bayes)

- For Gaussian naïve Bayes, you will arrive the following by taking the same derivation

