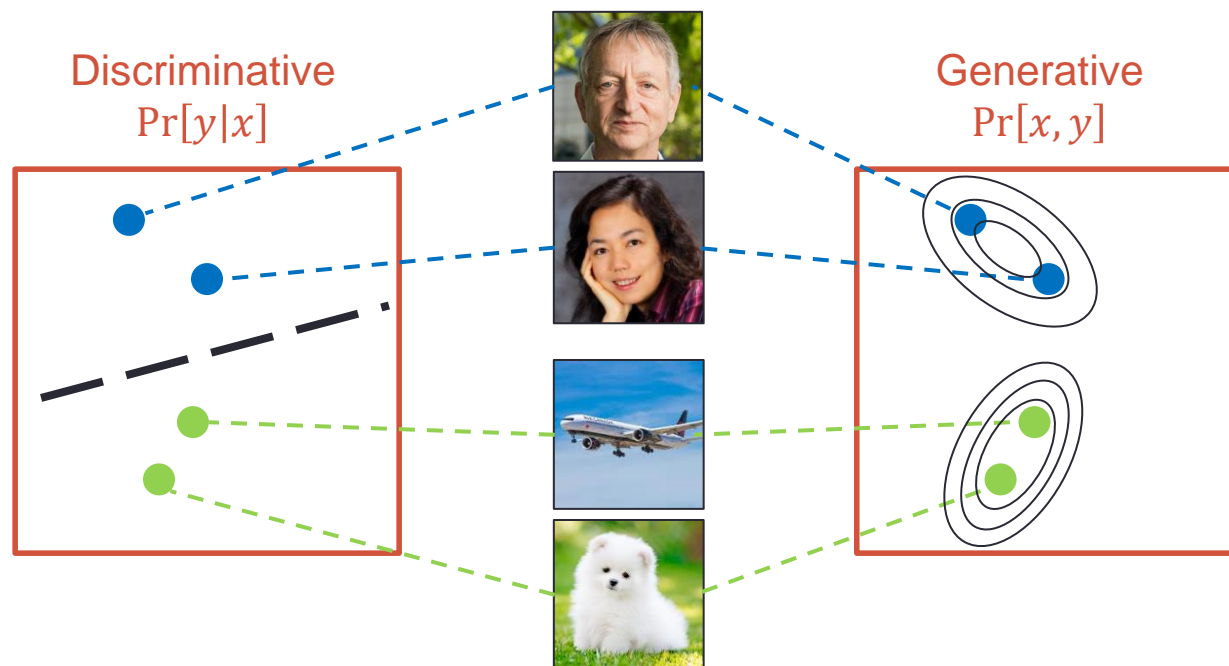


GENERATIVE MODELS

Generative models vs Discriminative models

- Generative models solve a harder task
- But sometimes it's easy to learn their parameters
- Can gain a better “understanding” of structures in data



Tutorial outline

- Maximum Likelihood Estimation (MLE)
- Naïve Bayes
- Maximum a Posteriori estimation (MAP)

MAXIMUM-LIKELIHOOD ESTIMATION

Maximum Likelihood Estimation (MLE)

- **Setting:** Independent draws $x_1, \dots, x_m \sim \mathcal{D}_\Theta$ from a parametric distribution
- **Likelihood:** Probability of the observed data given parameters

$$L(x_1, \dots, x_m | \Theta) = \Pr(S | \Theta) = \Pr(x_1, \dots, x_m | \Theta)$$

- **MLE** maximizes the likelihood

$$\hat{\Theta}_{\text{MLE}} = \operatorname{argmax}_{\Theta} L(x_1, \dots, x_m | \Theta)$$

- Equivalently,

$$\hat{\Theta}_{\text{MLE}} = \operatorname{argmax}_{\Theta} \ln L(x_1, \dots, x_m | \Theta)$$

Recall: MLE for Gaussian variables

- Given i.i.d Gaussian variables $x_1, \dots, x_m \sim \mathcal{N}(\mu, \sigma^2)$ we wish to estimate μ, σ^2
- Intuition:** given the observations 2, 9, 4 how would you “guess” μ ?

- In the lecture, we saw that:
 - The likelihood is:

$$L(S|\mu, \sigma^2) = \Pr(x_1, \dots, x_m|\mu, \sigma^2) = \prod_i \Pr(x_i|\mu, \sigma^2) = \prod_i \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{(x_i-\mu)^2}{2\sigma^2}\right\}$$

- The **log-likelihood** is: $\ln L(x_1, \dots, x_m|\mu, \sigma^2) = -m \cdot \ln \sigma\sqrt{2\pi} - \sum_i \frac{(x_i-\mu)^2}{2\sigma^2}$
- By simple differentiation, one can find the **MLEs**:

$$\hat{\mu}_{\text{MLE}} = \frac{1}{m} \sum_i x_i \quad , \quad \hat{\sigma}_{\text{MLE}}^2 = \frac{1}{m} \sum_i (x_i - \hat{\mu}_{\text{MLE}})^2$$

NAÏVE BAYES

Naïve Bayes

- Wish to find the **most probable label** by maximizing the posterior probability

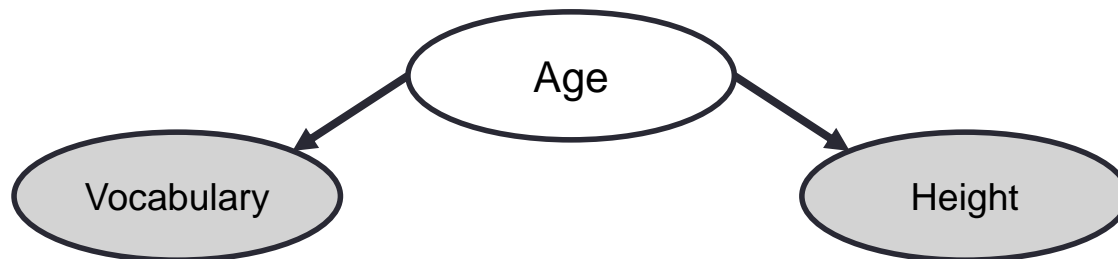
$$\hat{y} = h(\mathbf{x}) = \operatorname{argmax}_y \Pr(y|\mathbf{x}) \stackrel{\text{Bayes}}{=} \operatorname{argmax}_y \frac{\Pr(\mathbf{x}|y) \Pr(y)}{\Pr(\mathbf{x})} = \operatorname{argmax}_y \Pr(\mathbf{x}|y) \Pr(y)$$

- We make a **naïve assumption** – the coordinates are **conditionally independent**

$$\stackrel{\text{Naïve}}{=} \operatorname{argmax}_y \Pr(y) \prod_{k=1}^d \Pr(X[k] = x[k] | y)$$

Estimate these probabilities
(with MLE for instance)

- Often works well in practice, despite being a naïve assumption.
- Sometimes realistic! $\Pr(\text{height, vocabulary} | \text{age}) \approx \Pr(\text{height} | \text{age}) \cdot \Pr(\text{vocabulary} | \text{age})$



Demo: Iris dataset

- Three **classes**
 - Setosa, Versicolor, and Virginica
- Four **features**
 - Sepal length, sepal width, petal length, petal width (in cm)
 - Sepal is גביע,
 - Petal is כותרת



Iris Versicolor



Iris Setosa

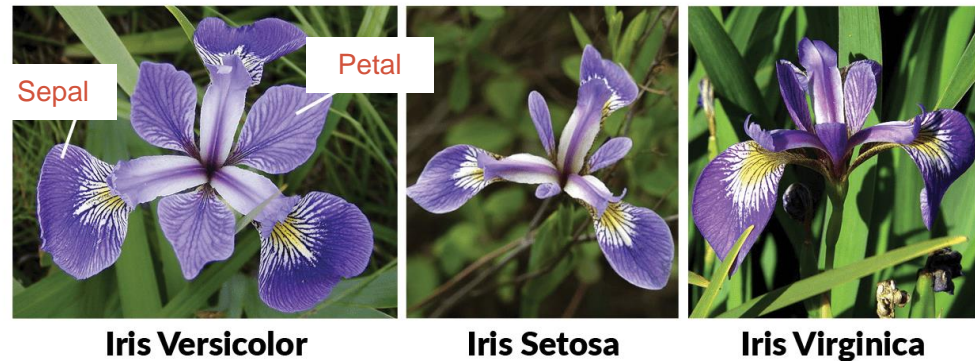


Iris Virginica

Source: [ML in R](#)

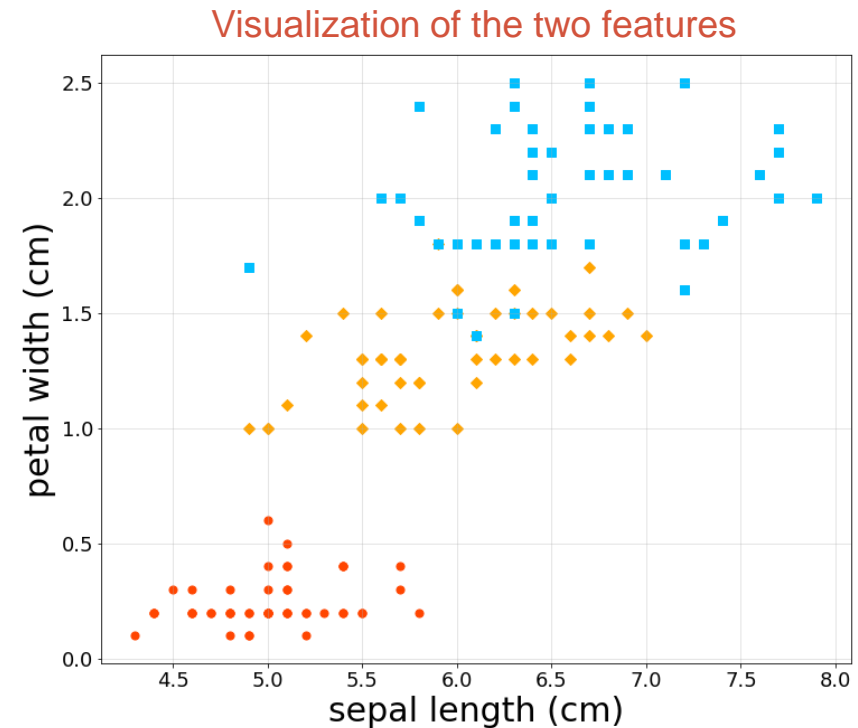
Demo: Iris dataset

- Three **classes**
 - Setosa, Versicolor, and Virginica



Source: [ML in R](#)

- Today we focus on two **features**
 - Sepal length ($x[1]$) and Petal width ($x[2]$)



Demo: Iris dataset

- Three **classes**
 - Setosa, Versicolor, and Virginica
- Today we focus on two **features**
 - Sepal length ($x[1]$) and Petal width ($x[2]$)
- **Naïve assumption:**
 - Given the species (y),
features are independent of each other



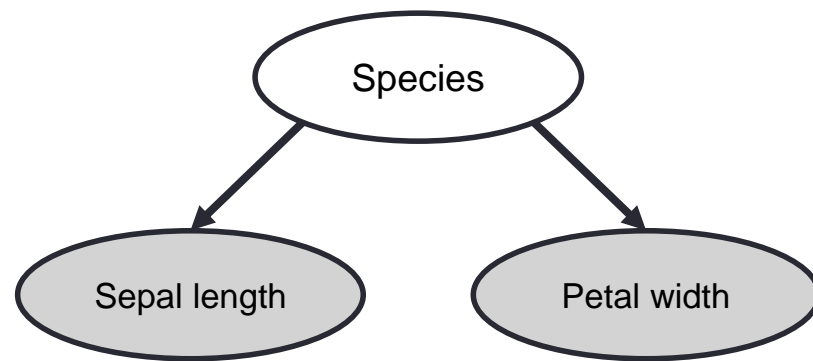
Iris Versicolor

Iris Setosa

Iris Virginica

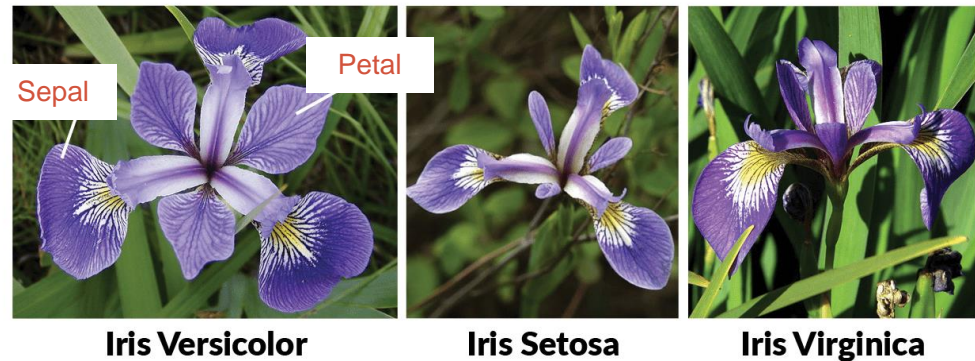
Source: ML in R

Graphical model:



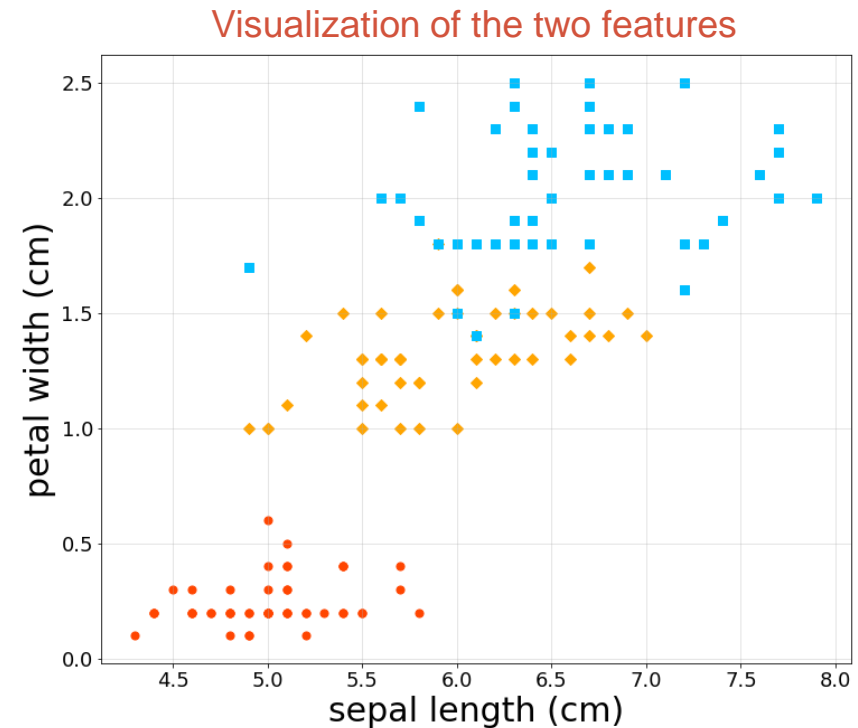
Demo: Iris dataset

- Three **classes**
 - Setosa, Versicolor, and Virginica



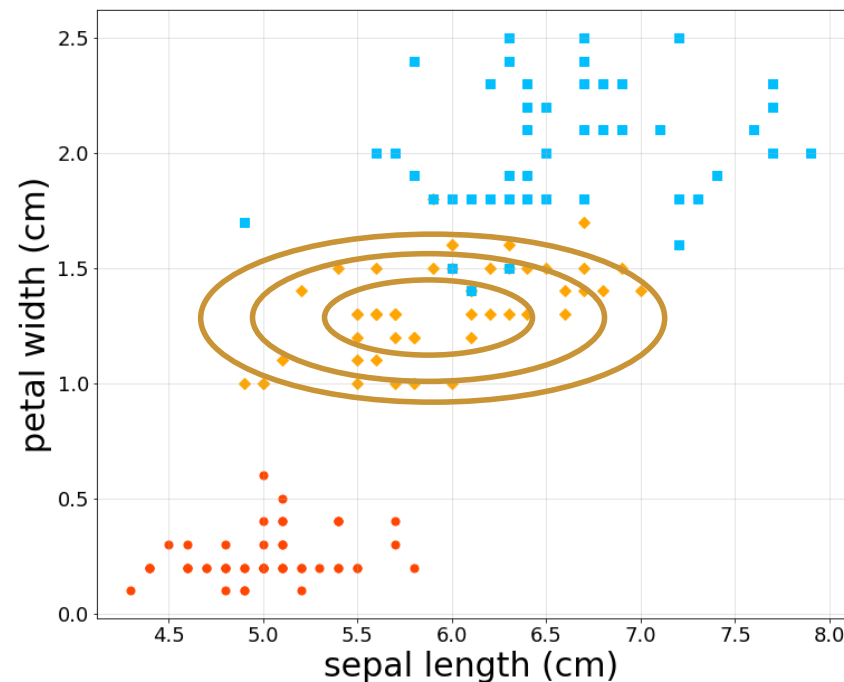
Source: [ML in R](#)

- Today we focus on two **features**
 - Sepal length ($x[1]$) and Petal width ($x[2]$)
- **Naïve assumption:**
 - Given the species (y),
features are independent of each other
- **Q:** in the plot, do they look independent?



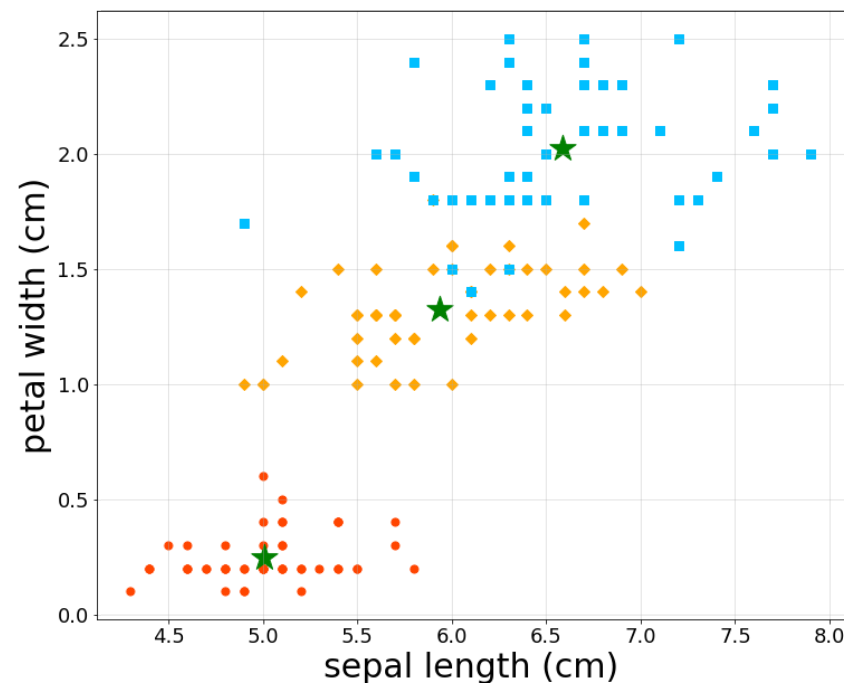
Demo: Gaussian Naïve Bayes

- **Naïve assumption:** given the species (y), features are independent
 - Use a Naïve Bayes classifier $\hat{y} = \operatorname{argmax}_y \Pr(y) \prod_{k=1}^d \Pr(X[k] = x[k] \mid y)$
- **Modeling:** assume probabilities $\Pr(X[k] = x[k] \mid y)$ are distributed $\mathcal{N}(\mu_y[k], \sigma_k^2)$
- **Goal:** fit multivariate Gaussians to the data
 - Estimate a different mean $\mu_y \in \mathbb{R}^d$ per class
 - Estimate a different variance $\sigma_k^2 \in \mathbb{R}_+$ per feature



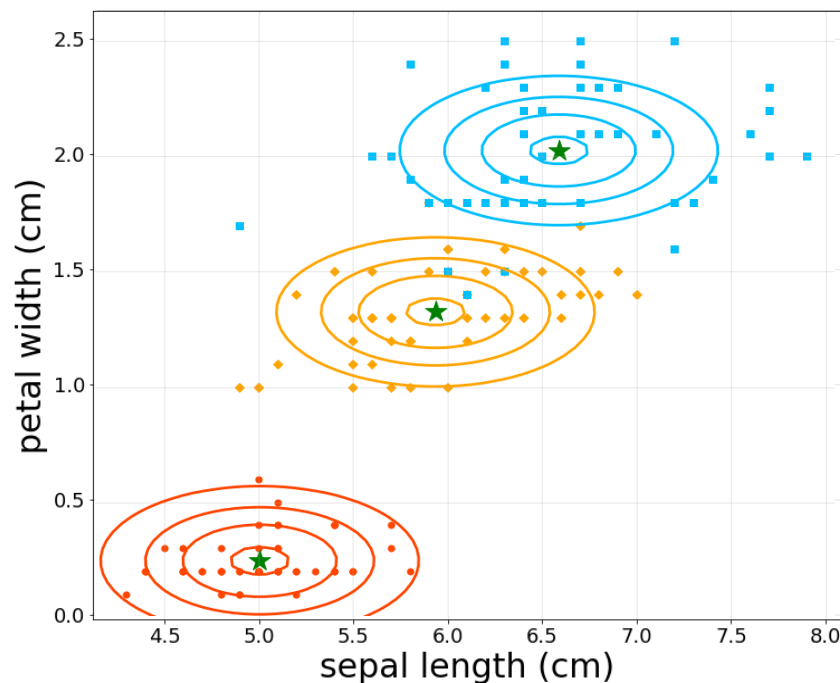
Demo: Gaussian Naïve Bayes

- **Naïve assumption:** given the species (y), features are independent
 - Use a Naïve Bayes classifier $\hat{y} = \operatorname{argmax}_y \Pr(y) \prod_{k=1}^d \Pr(X[k] = x[k] \mid y)$
- **Modeling:** assume probabilities $\Pr(X[k] = x[k] \mid y)$ are distributed $\mathcal{N}(\mu_y[k], \sigma_k^2)$
- **Goal:** fit multivariate Gaussians to the data
 - Estimate a different mean $\mu_y \in \mathbb{R}^d$ per class
 - Estimate a different variance $\sigma_k^2 \in \mathbb{R}_+$ per feature
- **Estimation:** maximize likelihood (MLE)
 - Means: $\hat{\mu}_y[k] = \frac{1}{\#\{y_i=y\}} \sum_{i: y_i=y} x_i[k]$



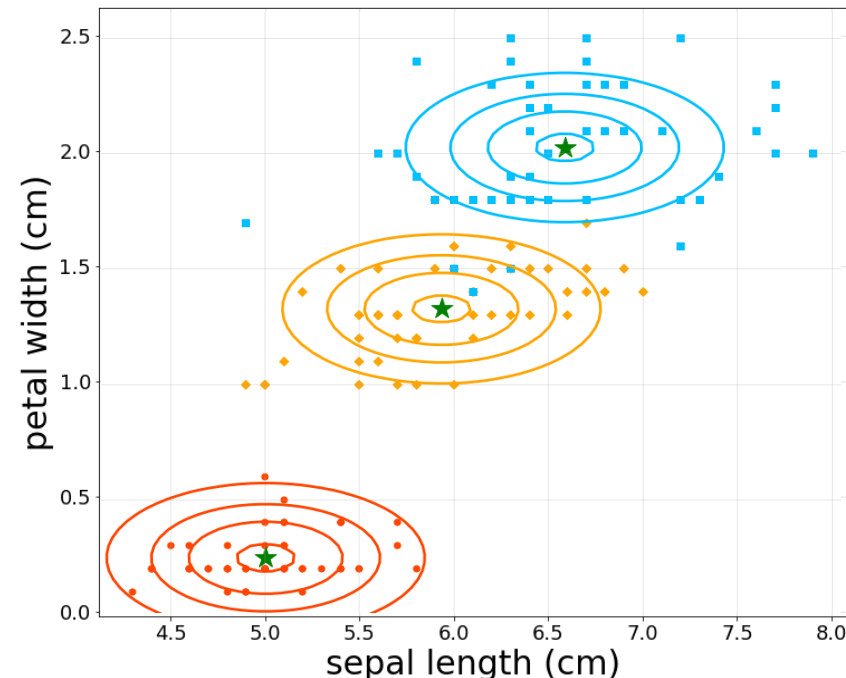
Demo: Gaussian Naïve Bayes

- **Naïve assumption:** given the species (y), features are independent
 - Use a Naïve Bayes classifier $\hat{y} = \operatorname{argmax}_y \Pr(y) \prod_{k=1}^d \Pr(X[k] = x[k] \mid y)$
- **Modeling:** assume probabilities $\Pr(X[k] = x[k] \mid y)$ are distributed $\mathcal{N}(\mu_y[k], \sigma_k^2)$
- **Goal:** fit multivariate Gaussians to the data
 - Estimate a different mean $\mu_y \in \mathbb{R}^d$ per class
 - Estimate a different variance $\sigma_k^2 \in \mathbb{R}_+$ per feature
- **Estimation:** maximize likelihood (MLE)
 - Means: $\hat{\mu}_y[k] = \frac{1}{\#\{y_i=y\}} \sum_{i: y_i=y} x_i[k]$
 - Variances: $\hat{\sigma}_k^2 = \frac{1}{m} \sum_i (x_i[k] - \hat{\mu}_{y_i}[k])^2$ (extra)
- **Q:** why are the Gaussians “axis-aligned”?



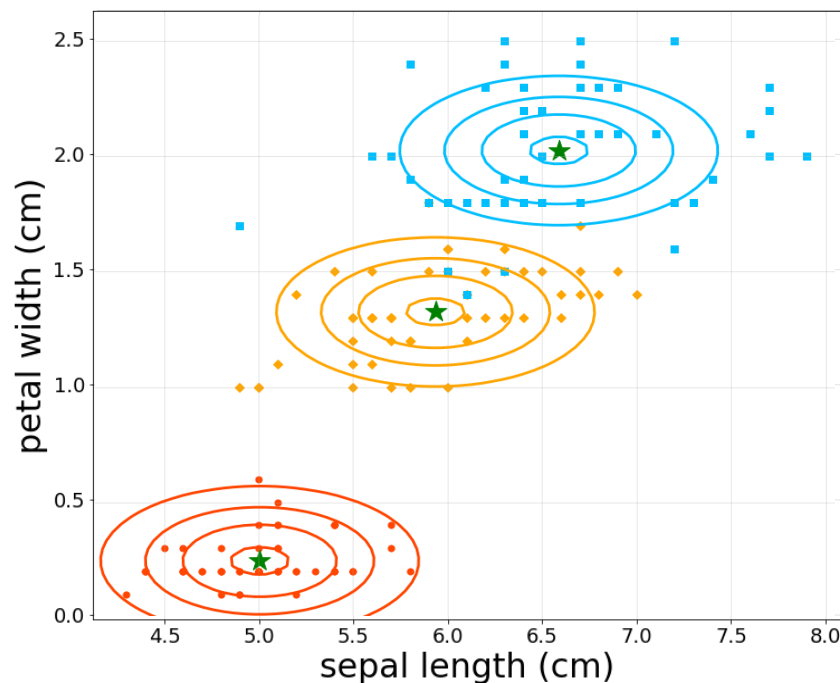
Demo: Gaussian Naïve Bayes

- **Naïve assumption:** given the species (y), features are independent
 - Use a Naïve Bayes classifier $\hat{y} = \operatorname{argmax}_y \Pr(y) \prod_{k=1}^d \Pr(X[k] = x[k] \mid y)$
- **Modeling:** assume probabilities $\Pr(X[k] = x[k] \mid y)$ are distributed $\mathcal{N}(\mu_y[k], \sigma_k^2)$
- **Goal:** fit multivariate Gaussians to the data
 - Estimate a different mean $\mu_y \in \mathbb{R}^d$ per class
 - Estimate a different variance $\sigma_k^2 \in \mathbb{R}_+$ per feature
- **Estimation:** maximize likelihood (MLE)
 - Means: $\hat{\mu}_y[k] = \frac{1}{\#\{y_i=y\}} \sum_{i: y_i=y} x_i[k]$
 - Variances: $\hat{\sigma}_k^2 = \frac{1}{m} \sum_i (x_i[k] - \hat{\mu}_{y_i}[k])^2$ (extra)
- **Q:** why all Gaussians “look” the same?



Demo: Gaussian Naïve Bayes

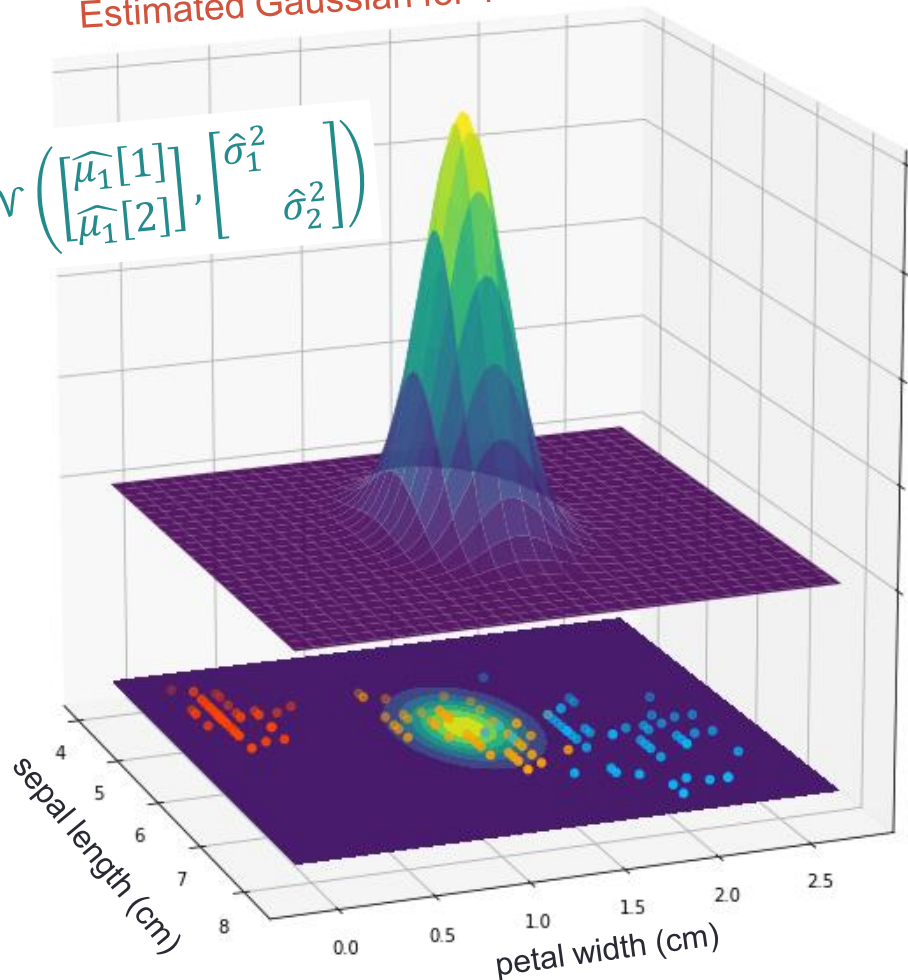
- **Naïve assumption:** given the species (y), features are independent
 - Use a Naïve Bayes classifier $\hat{y} = \operatorname{argmax}_y \Pr(y) \prod_{k=1}^d \Pr(X[k] = x[k] \mid y)$
- **Modeling:** assume probabilities $\Pr(X[k] = x[k] \mid y)$ are distributed $\mathcal{N}(\mu_y[k], \sigma_k^2)$
- **Goal:** fit multivariate Gaussians to the data
 - Estimate a different mean $\mu_y \in \mathbb{R}^d$ per class
 - Estimate a different variance $\sigma_k^2 \in \mathbb{R}_+$ per feature
- **Estimation:** maximize likelihood (MLE)
 - Means: $\widehat{\mu}_y[k] = \frac{1}{\#\{y_i=y\}} \sum_{i: y_i=y} x_i[k]$
 - Variances: $\widehat{\sigma}_k^2 = \frac{1}{m} \sum_i (x_i[k] - \widehat{\mu}_{y_i}[k])^2$ (extra)
 - Marginal: $\widehat{\Pr}(y) = \frac{1}{m} \#\{y_i = y\}$



Demo: Estimated Gaussians

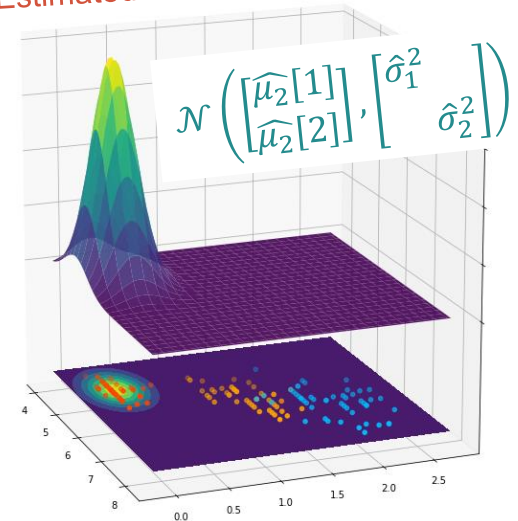
Estimated Gaussian for 1st class

$$\mathcal{N}\left(\begin{bmatrix} \hat{\mu}_1[1] \\ \hat{\mu}_1[2] \end{bmatrix}, \begin{bmatrix} \hat{\sigma}_1^2 & \\ & \hat{\sigma}_2^2 \end{bmatrix}\right)$$



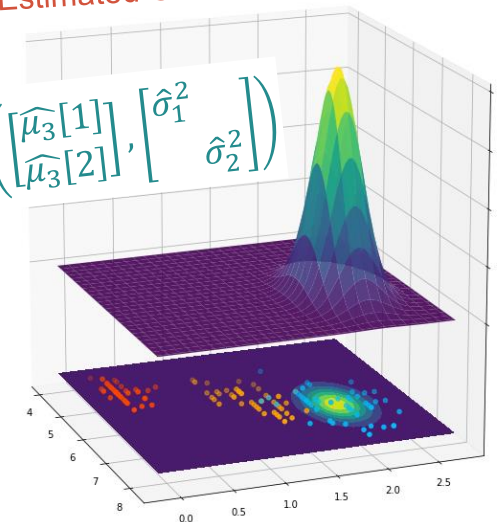
Estimated Gaussian for 2nd class

$$\mathcal{N}\left(\begin{bmatrix} \hat{\mu}_2[1] \\ \hat{\mu}_2[2] \end{bmatrix}, \begin{bmatrix} \hat{\sigma}_1^2 & \\ & \hat{\sigma}_2^2 \end{bmatrix}\right)$$



Estimated Gaussian for 3rd class

$$\mathcal{N}\left(\begin{bmatrix} \hat{\mu}_3[1] \\ \hat{\mu}_3[2] \end{bmatrix}, \begin{bmatrix} \hat{\sigma}_1^2 & \\ & \hat{\sigma}_2^2 \end{bmatrix}\right)$$

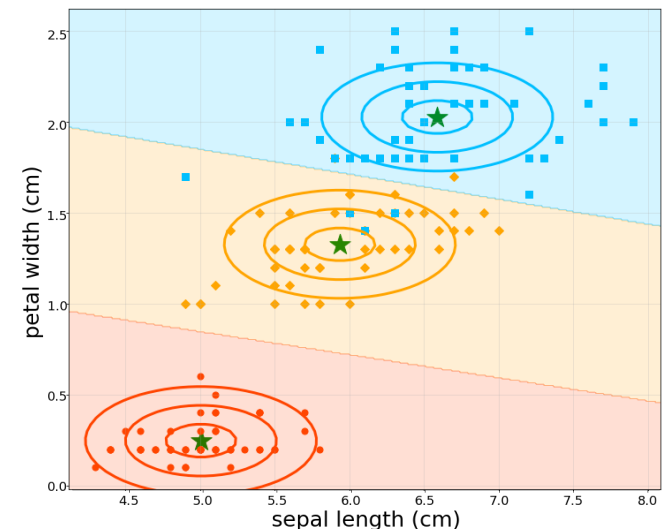


Demo: Making a prediction

- Naïve Bayes rule: $\hat{y} = \operatorname{argmax}_y \Pr(y|\mathbf{x}) = \operatorname{argmax}_y \Pr(y) \prod_{k=1}^d \Pr(X[k] = x[k] | y)$
- **Prediction** using our estimators:

$$h(\mathbf{x}) = \operatorname{argmax}_y \widehat{\Pr}(y) \prod_{k=1}^d \frac{1}{\widehat{\sigma}_k \sqrt{2\pi}} \exp \left\{ -\frac{(\mathbf{x} - \widehat{\mu}_y)^2}{2\widehat{\sigma}_k^2} \right\}$$

- The predictor asks: which Gaussian gives the maximal probability to seeing \mathbf{x} ?
(normalized by the “prior”/marginal probability)
- Assuming same covariance for all classes,
decision boundaries are **linear** (proof in lecture)
- **Q:** What is the training error here?



MAXIMUM A POSTERIORI ESTIMATION

Recall: Least squares as MLE

- We saw the following theorem (lecture 09):

- Assuming a noisy linear model:

$$y_i = \mathbf{x}_i^\top \mathbf{w} + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, 1) \text{ (i.i.d)}$$

- Note: a sample \mathbf{x}_i is not random but the noise ε_i is.

- Solving Least Squares (LS) is equivalent to Maximum-Likelihood Estimation (MLE):

$$\hat{\mathbf{w}}_{\text{LS}} \triangleq \operatorname{argmin}_{\mathbf{w}} \frac{1}{m} \sum_i (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 = \operatorname{argmax}_{\mathbf{w}} L(\mathbf{w}; S) \triangleq \hat{\mathbf{w}}_{\text{MLE}}$$

- Exercise (like in lecture 09):

- Prove that the likelihood is as follows, by justifying the equalities below:

$$\begin{aligned} L(\mathbf{w}; S) &\triangleq P(\{(\mathbf{x}_i, y_i)\}_{i=1}^m | \mathbf{w}) = \prod_{i=1}^m P(y_i, | \mathbf{x}_i, \mathbf{w}) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i^\top \mathbf{w} - y_i)^2 \right\} \\ &= (2\pi)^{-\frac{m}{2}} \cdot \exp \left\{ -\frac{1}{2} \sum_{i=1}^m (\mathbf{x}_i^\top \mathbf{w} - y_i)^2 \right\} \end{aligned}$$

Recall: Least squares as MLE

- We saw the following theorem (lecture 09):

- Assuming a noisy linear model:

$$y_i = \mathbf{x}_i^\top \mathbf{w} + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, 1) \text{ (i.i.d.)}$$

- Note: a sample \mathbf{x}_i is not random but the noise ε_i is.

- Solving Least Squares (LS) is equivalent to Maximum-Likelihood Estimation (MLE):

$$\hat{\mathbf{w}}_{\text{LS}} \triangleq \operatorname{argmin}_{\mathbf{w}} \frac{1}{m} \sum_i (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 = \operatorname{argmax}_{\mathbf{w}} L(\mathbf{w}; S) \triangleq \hat{\mathbf{w}}_{\text{MLE}}$$

- Exercise (like in lecture 09):

- The likelihood is:

$$L(\mathbf{w}; S) = (2\pi)^{-\frac{m}{2}} \cdot \exp \left\{ -\frac{1}{2} \sum_{i=1}^m (\mathbf{x}_i^\top \mathbf{w} - y_i)^2 \right\}$$

- Prove the theorem above.

$$\hat{\mathbf{w}}_{\text{MLE}} \triangleq \operatorname{argmax}_{\mathbf{w}} L(\mathbf{w}; S) =$$

$$= \operatorname{argmin}_{\mathbf{w}} \frac{1}{m} \sum_i (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 \triangleq \hat{\mathbf{w}}_{\text{LS}}$$

Recall: Least squares as MLE

- We saw the following theorem (lecture 09):

- Assuming a noisy linear model:

$$y_i = \mathbf{x}_i^\top \mathbf{w} + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, 1) \text{ (i.i.d)}$$

- Note: a sample \mathbf{x}_i is not random but the noise ε_i is.

- Solving Least Squares (LS) is equivalent to Maximum-Likelihood Estimation (MLE):

$$\hat{\mathbf{w}}_{\text{LS}} \triangleq \operatorname{argmin}_{\mathbf{w}} \frac{1}{m} \sum_i (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 = \operatorname{argmax}_{\mathbf{w}} L(\mathbf{w}; S) \triangleq \hat{\mathbf{w}}_{\text{MLE}}$$

- Exercise (like in lecture 09):

- The likelihood is:
$$L(\mathbf{w}; S) = (2\pi)^{-\frac{m}{2}} \cdot \exp \left\{ -\frac{1}{2} \sum_{i=1}^m (\mathbf{x}_i^\top \mathbf{w} - y_i)^2 \right\}$$

- Prove the theorem above.

$$\begin{aligned} \hat{\mathbf{w}}_{\text{MLE}} &\triangleq \operatorname{argmax}_{\mathbf{w}} L(\mathbf{w}; S) = \operatorname{argmax}_{\mathbf{w}} \ln \left[(2\pi)^{-\frac{m}{2}} \cdot \exp \left\{ -\frac{1}{2} \sum_{i=1}^m (\mathbf{x}_i^\top \mathbf{w} - y_i)^2 \right\} \right] \\ &= \operatorname{argmax}_{\mathbf{w}} \left(-\frac{m}{2} \ln(2\pi) - \frac{1}{2} \sum_{i=1}^m (\mathbf{x}_i^\top \mathbf{w} - y_i)^2 \right) = \operatorname{argmin}_{\mathbf{w}} \frac{1}{m} \sum_i (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 \triangleq \hat{\mathbf{w}}_{\text{LS}} \end{aligned}$$

Maximum a Posteriori (MAP) Estimation

- Maximizes the **posterior probability** of the parameters given the observations

$$\hat{\Theta}_{\text{MAP}} = \operatorname{argmax}_{\Theta} \Pr[\Theta|S] = \operatorname{argmax}_{\Theta} \frac{\Pr[S|\Theta] \Pr[\Theta]}{\Pr[S]} = \operatorname{argmax}_{\Theta} \underbrace{\Pr[S|\Theta]}_{\triangleq L(\Theta; S)} \Pr[\Theta]$$

- Assumes a prior on the parameters themselves!
- Notice the difference from **MLE**

$$\hat{\Theta}_{\text{MLE}} = \operatorname{argmax}_{\Theta} \Pr[S|\Theta]$$

- **Q:** When are they equivalent?

Exercise: Ridge regression as MAP

- Like before, we assume a noisy linear model: $y_i = \mathbf{x}_i^\top \mathbf{w} + \varepsilon_i$, $\varepsilon_i \sim \mathcal{N}(0, 1)$ (i.i.d)
 - Recall, The likelihood is:
$$L(\mathbf{w}; S) = (2\pi)^{-\frac{m}{2}} \cdot \exp \left\{ -\frac{1}{2} \sum_{i=1}^m (\mathbf{x}_i^\top \mathbf{w} - y_i)^2 \right\}$$
- We further assume a prior on the weights: $w[k] \sim \mathcal{N}(0, 1/\lambda_m)$, $\lambda > 0$.

1. Express the prior PDF $p(\mathbf{w})$

Exercise: Ridge regression as MAP

- Like before, we assume a noisy linear model: $y_i = \mathbf{x}_i^\top \mathbf{w} + \varepsilon_i$, $\varepsilon_i \sim \mathcal{N}(0, 1)$ (i.i.d)
 - Recall, The likelihood is:
$$L(\mathbf{w}; S) = (2\pi)^{-\frac{m}{2}} \cdot \exp \left\{ -\frac{1}{2} \sum_{i=1}^m (\mathbf{x}_i^\top \mathbf{w} - y_i)^2 \right\}$$
- We further assume a prior on the weights: $w[k] \sim \mathcal{N}(0, 1/\lambda_m)$, $\lambda > 0$.

1. Express the prior PDF $p(\mathbf{w}) = \prod_{k=1}^d p(w[k])$

$$\begin{aligned}
 &= \prod_{k=1}^d \sqrt{\frac{\lambda m}{2\pi}} \exp \left\{ -\frac{\lambda}{2} w[k]^2 \right\} \\
 &= \left(\frac{\lambda m}{2\pi} \right)^{\frac{d}{2}} \exp \left\{ -\frac{\lambda m}{2} \sum_{k=1}^d w[k]^2 \right\}
 \end{aligned}$$

Exercise: Ridge regression as MAP

- Like before, we assume a noisy linear model: $y_i = \mathbf{x}_i^\top \mathbf{w} + \varepsilon_i$, $\varepsilon_i \sim \mathcal{N}(0, 1)$ (i.i.d)

- Recall, The likelihood is:
$$L(\mathbf{w}; S) = (2\pi)^{-\frac{m}{2}} \cdot \exp \left\{ -\frac{1}{2} \sum_{i=1}^m (\mathbf{x}_i^\top \mathbf{w} - y_i)^2 \right\}$$

- We further assume a prior on the weights: $w[k] \sim \mathcal{N}(0, 1/\lambda_m)$, $\lambda > 0$.

- The prior PDF is $p(\mathbf{w}) = \left(\frac{\lambda m}{2\pi}\right)^{\frac{d}{2}} \exp \left\{ -\frac{\lambda m}{2} \sum_{k=1}^d w[k]^2 \right\}$

- Prove:** under the assumptions above,

solving Ridge regression is equivalent to Maximum a posteriori Estimation (MAP):

$$\hat{\mathbf{w}}_{\text{Ridge}} \triangleq \operatorname{argmin}_{\mathbf{w}} \frac{1}{m} \sum_i (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|_2^2 = \operatorname{argmax}_{\mathbf{w}} L(\mathbf{w}; S) p(\mathbf{w}) \triangleq \hat{\mathbf{w}}_{\text{MAP}}$$

Exercise: Ridge regression as MAP

- Like before, we assume a noisy linear model: $y_i = \mathbf{x}_i^\top \mathbf{w} + \varepsilon_i$, $\varepsilon_i \sim \mathcal{N}(0, 1)$ (i.i.d)
 - Recall, The likelihood is:
$$L(\mathbf{w}; S) = (2\pi)^{-\frac{m}{2}} \cdot \exp \left\{ -\frac{1}{2} \sum_{i=1}^m (\mathbf{x}_i^\top \mathbf{w} - y_i)^2 \right\}$$
- We further assume a prior on the weights: $w[k] \sim \mathcal{N}(0, 1/\lambda m)$, $\lambda > 0$.
 - The prior PDF is $p(\mathbf{w}) = \left(\frac{\lambda m}{2\pi} \right)^{\frac{d}{2}} \exp \left\{ -\frac{\lambda m}{2} \sum_{k=1}^d w[k]^2 \right\}$
 - Proof:** $\hat{\mathbf{w}}_{\text{MAP}} \triangleq \operatorname{argmax}_{\mathbf{w}} L(\mathbf{w}; S) p(\mathbf{w})$

$$= \operatorname{argmax}_{\mathbf{w}} (\ln L(\mathbf{w}; S) + \ln p(\mathbf{w}))$$

$$= \operatorname{argmax}_{\mathbf{w}} \left(-\frac{m}{2} \ln(2\pi) - \frac{1}{2} \sum_{i=1}^m (\mathbf{x}_i^\top \mathbf{w} - y_i)^2 + \frac{d}{2} \ln \left(\frac{\lambda m}{2\pi} \right) - \frac{\lambda m}{2} \sum_{k=1}^d w[k]^2 \right)$$

$$= \operatorname{argmin}_{\mathbf{w}} \frac{1}{m} \sum_i (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|_2^2 \triangleq \hat{\mathbf{w}}_{\text{Ridge}}$$

Summary

- **MLE** finds the most likely model to have generated the data

$$\hat{\Theta}_{\text{MLE}} = \operatorname{argmax}_{\Theta} \Pr[S|\Theta]$$

- **MAP** finds the most probable model based on the data

$$\hat{\Theta}_{\text{MAP}} = \operatorname{argmax}_{\Theta} \Pr[\Theta|S] = \operatorname{argmax}_{\Theta} \Pr[S|\Theta] \Pr[\Theta]$$

- **Naïve Bayes** makes a naïve assumption that the features are conditionally independent given the label

$$\Pr[\mathbf{x}|y] = \prod_{k=1}^d \Pr[X[k] = x_k | Y = y]$$