

Classifiers

Readings for today

- Chapter 4: Classification. James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). An introduction to statistical learning: with applications in R (Vol. 6). New York: Springer.

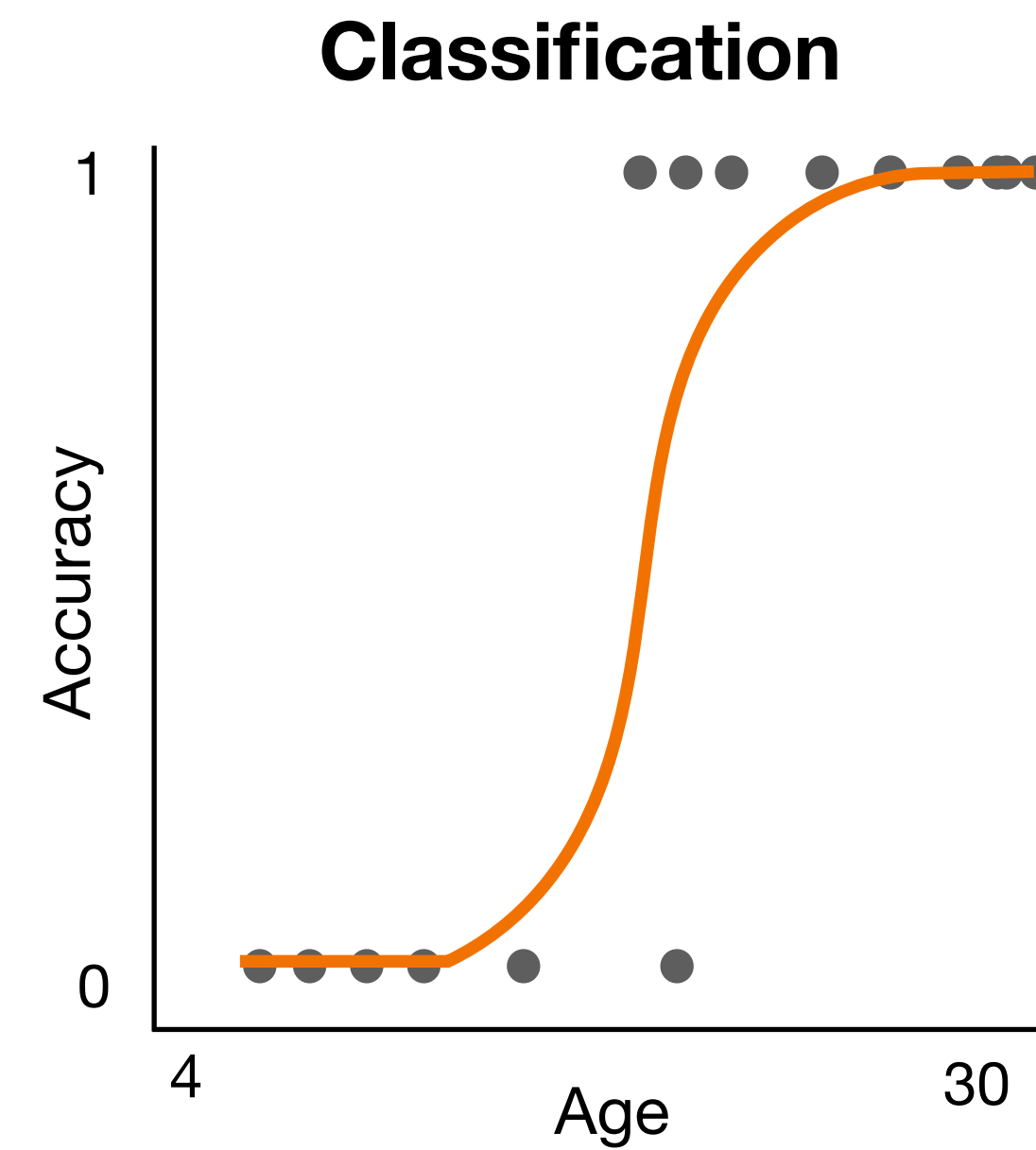
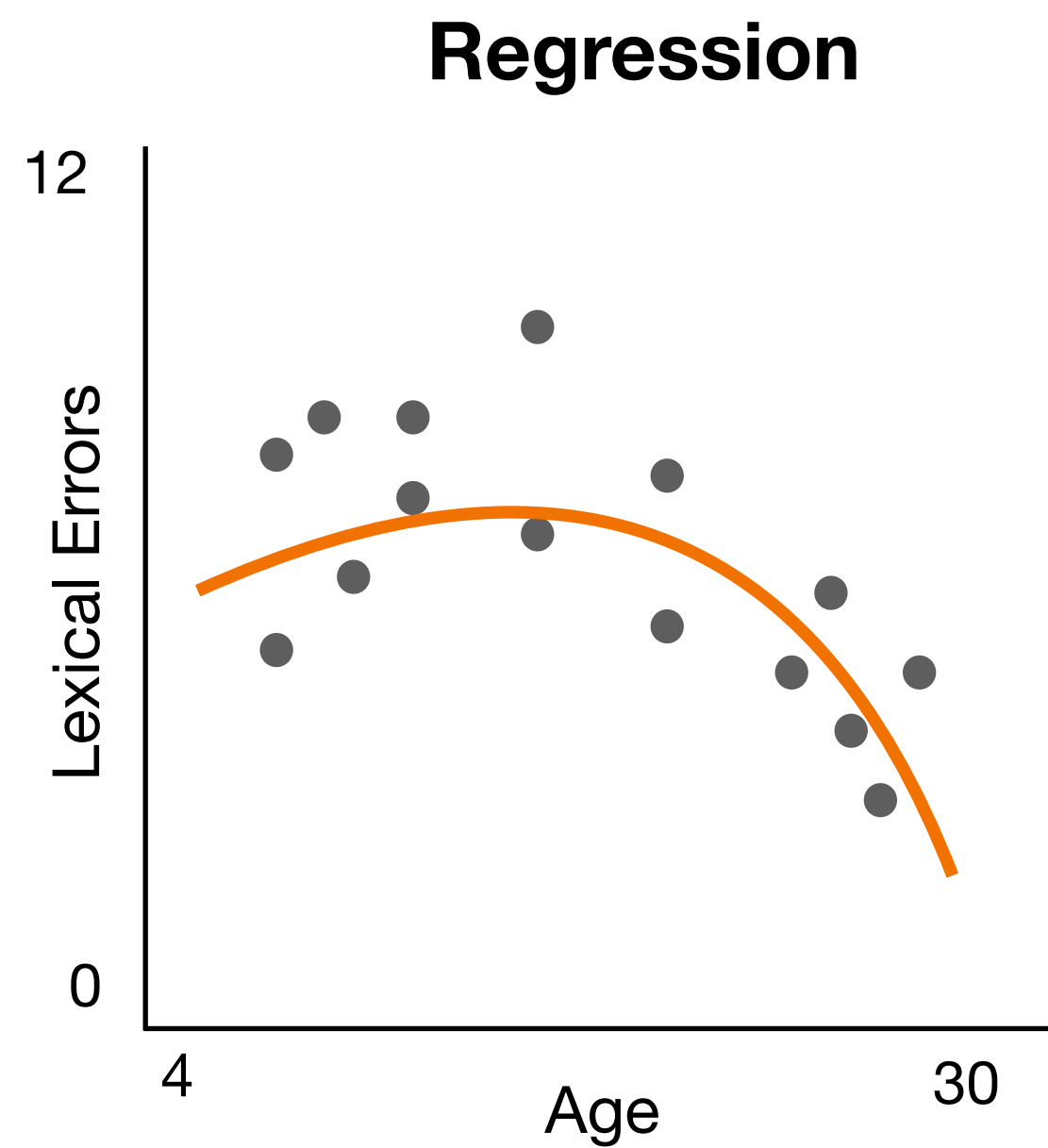
Topics

1. Logistic regression
2. Linear & quadratic discriminant analysis

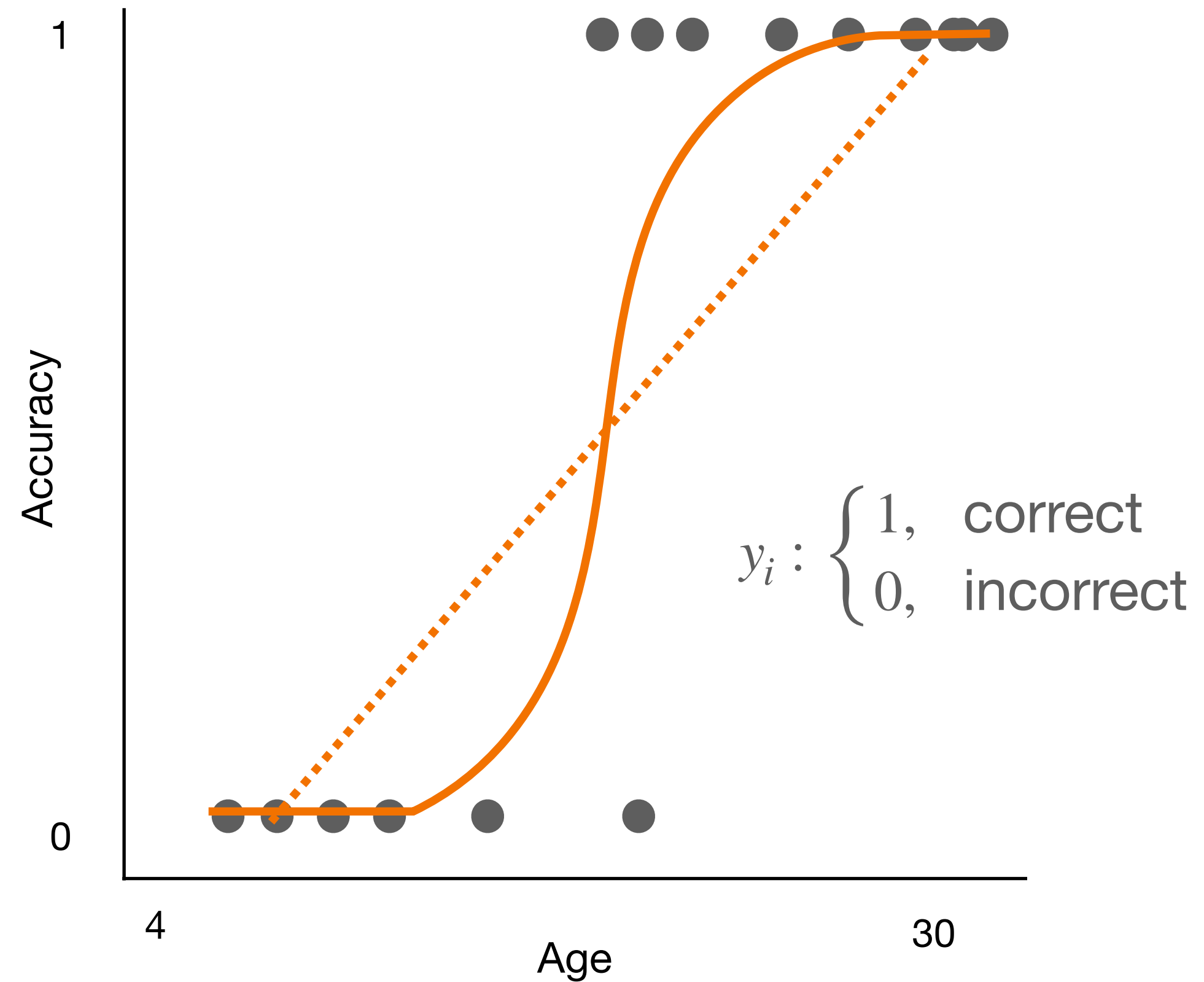
Logistic regression

Classes of models

		Y	
		Quantitative	Qualitative
Regression		✓	
Classification			✓



Logistic regression



Linear regression

$$f(y_i | x_i, \beta_1, \beta_0, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{\sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2}{2\sigma^2}}$$

Residual errors are normally distributed

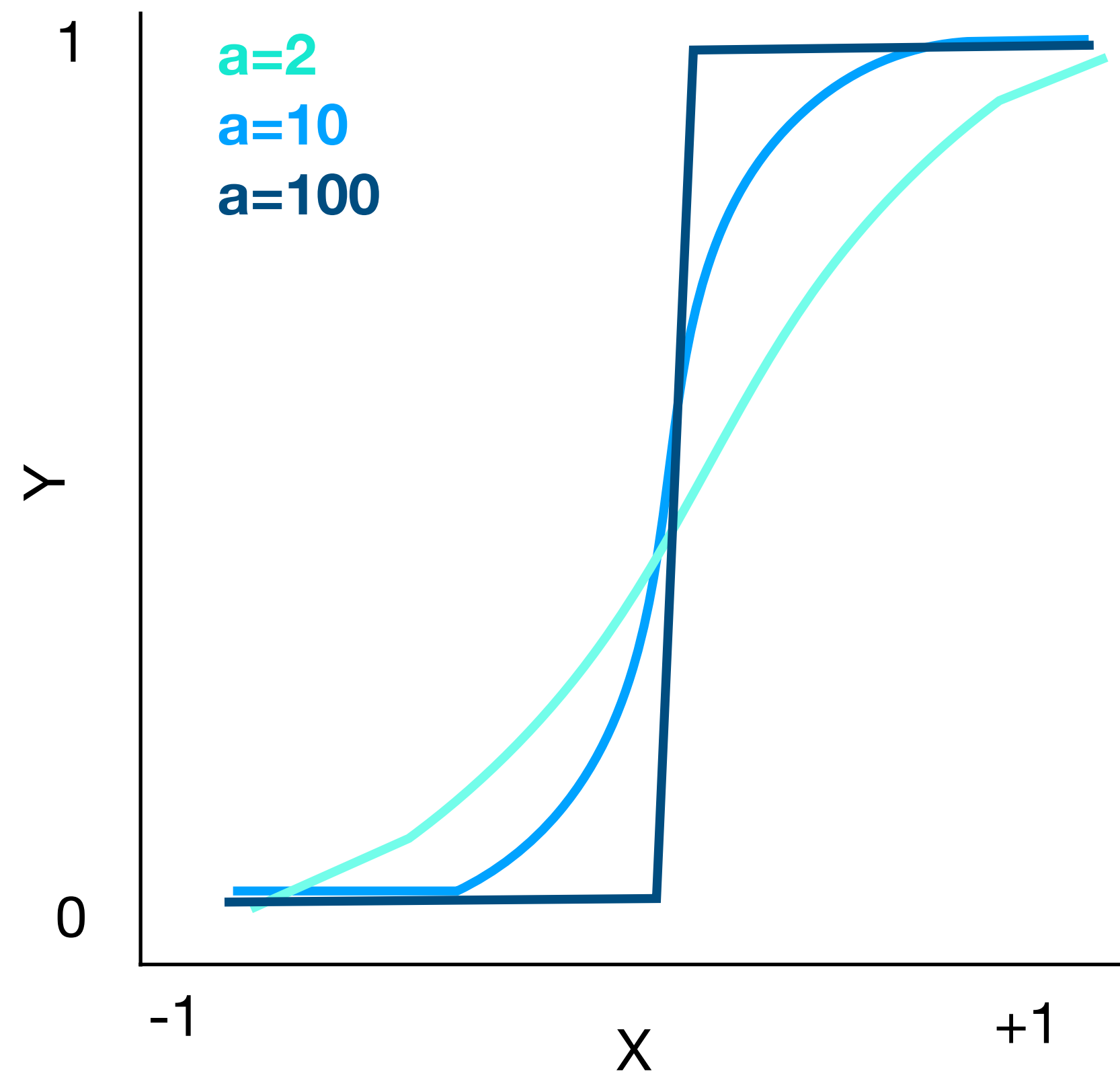
Logistic regression

$$p(y = 1 | x) = f(y_i | x_i, \beta_1, \beta_0, \sigma) = \frac{e^{\hat{\beta}_0 + \hat{\beta}_1 X}}{1 + e^{\hat{\beta}_0 + \hat{\beta}_1 X}}$$

Residual errors are binomial (i.e., exist in one of 2 states)

The standard logistic function

Standard logistic function: $p(y = 1 | x) = \frac{e^{ax}}{1 + e^{ax}}$

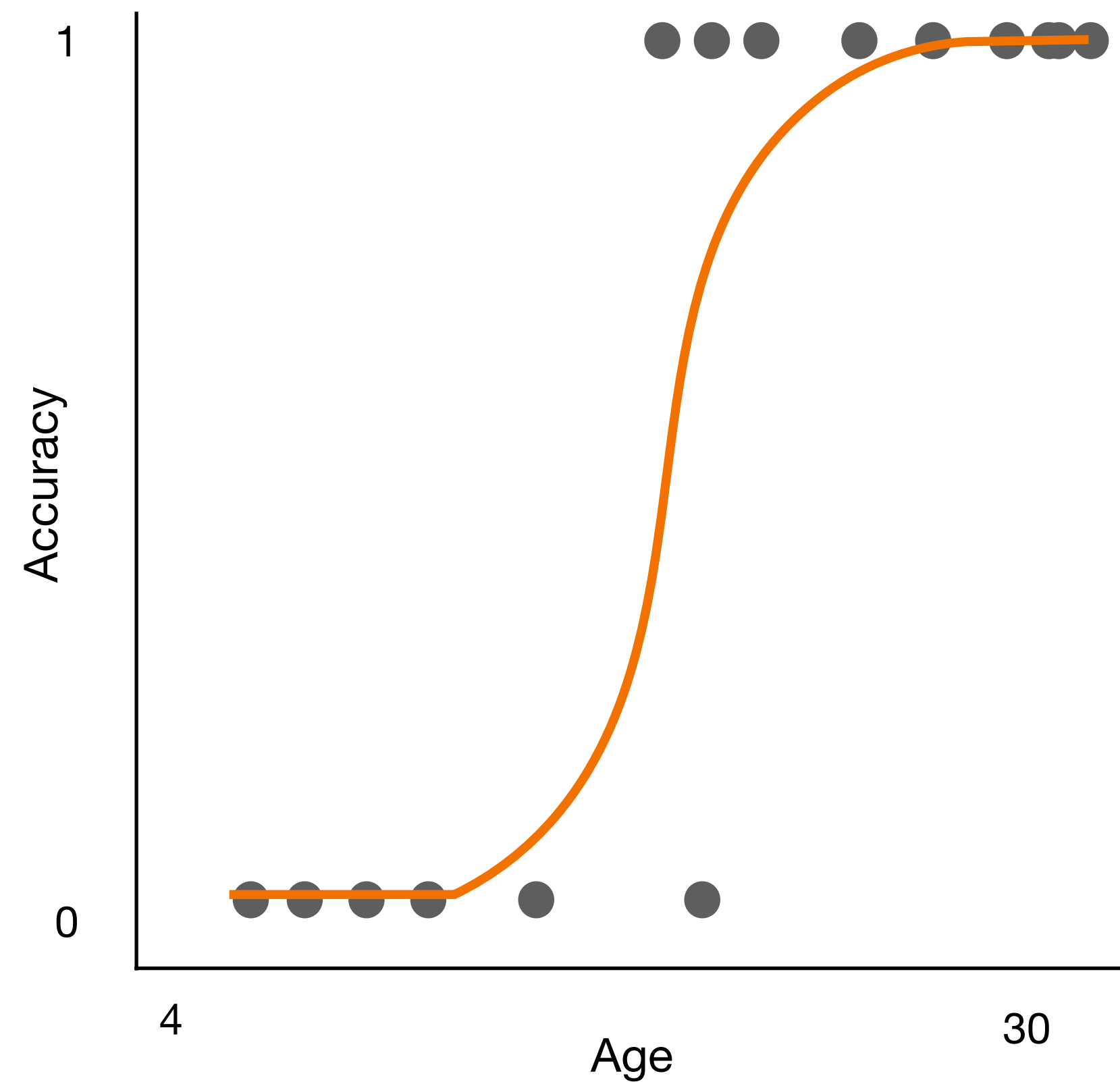


- Lower *growth rate* parameter (a), means slower transition between states of y (i.e., greater uncertainty).
- In a regression context better separation of y as x changes leads to larger a (i.e., $a = \sum_{j=1}^p \hat{\beta}_j X_j + \hat{\beta}_0$)

Odds

$$odds(y = 1) = \frac{p(y = 1)}{1 - p(y = 1)} = e^{ax}$$
$$\ln(odds(y = 1)) = \ln\left(\frac{p(y = 1)}{1 - p(y = 1)}\right) = ax = \hat{\beta}_0 + \sum_{j=1}^p \hat{\beta}_j X_j$$

Logistic regression



Assumptions:

1. $f(X)$ describes a linear relationship between X and Y .
2. Y is binary/dichotomous.
3. There is no collinearity between features in X .
4. $f(X)$ is stationary.

Interpretation

Linear regression

$$\hat{y}_i = \beta_0 + \sum_{j=1}^p \beta_j x_{i,j} = 1 + 0.5x_i$$

Interpretation

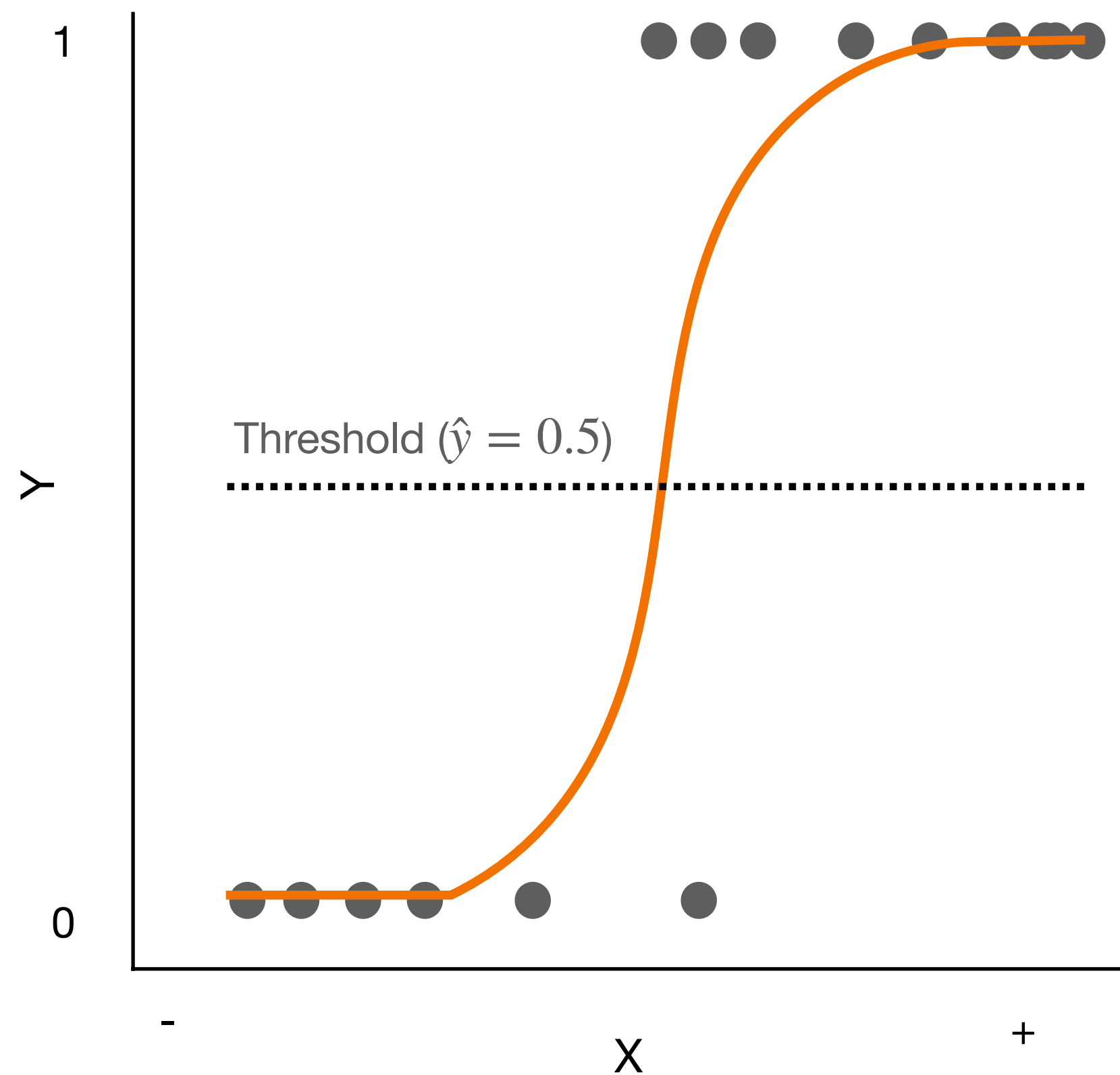
For a 1 unit change in x , y changes 0.5, added to a baseline of 1.

Logistic regression

$$\hat{y}_i = e^{\hat{\beta}_0 + \sum_{j=1}^p \beta_j x_{i,j}} = e^{1+0.5x_i}$$

For a 1 unit change in x , the log odds of $y = 1$ changes by 0.5, multiplied by a baseline odds of e^1

Prediction

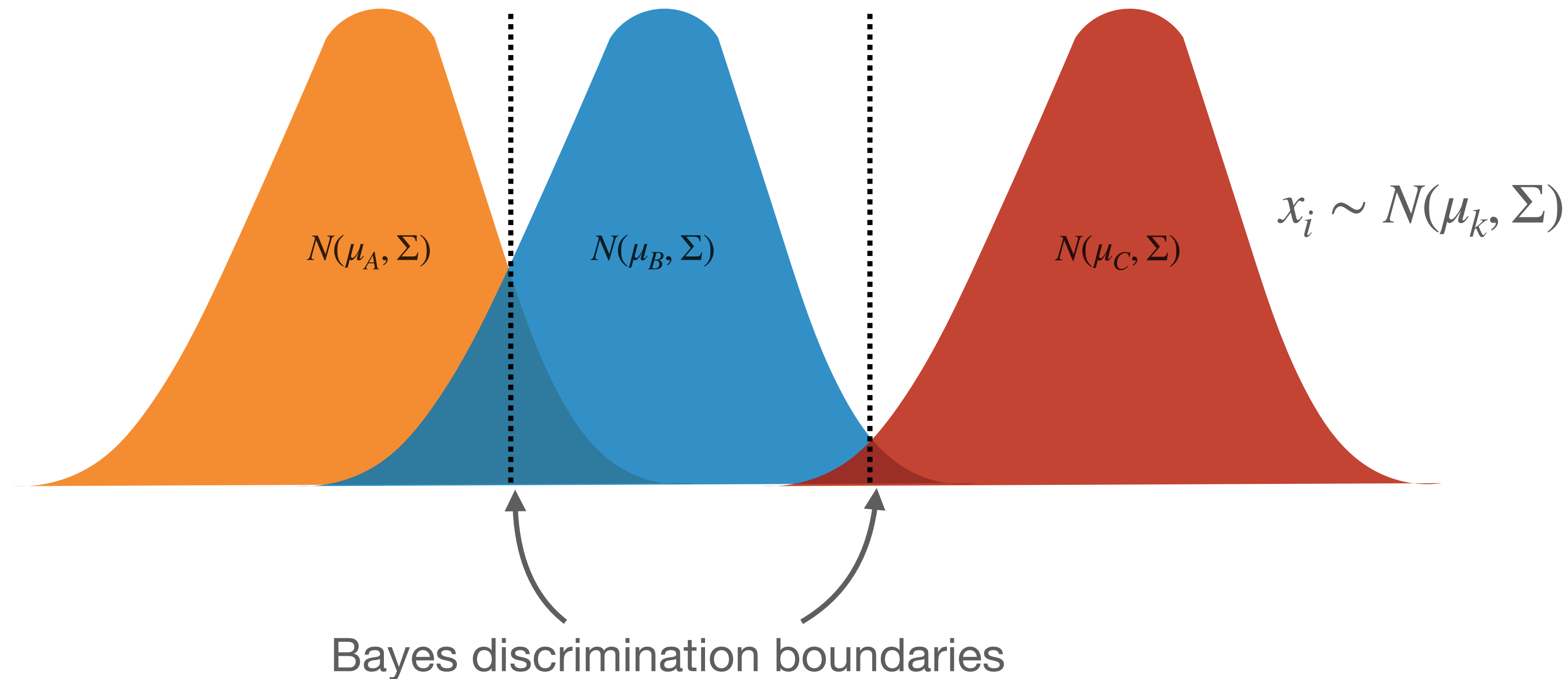


- $\hat{f}(X)$ is still a continuous function.
- In order to predict y , an *a priori* threshold needs to be determined to force a category for every $\hat{f}(x_i) = \hat{y}_i$.
- Position of this threshold determines the bias of your predictions.

Linear & Quadratic Discriminant Analysis

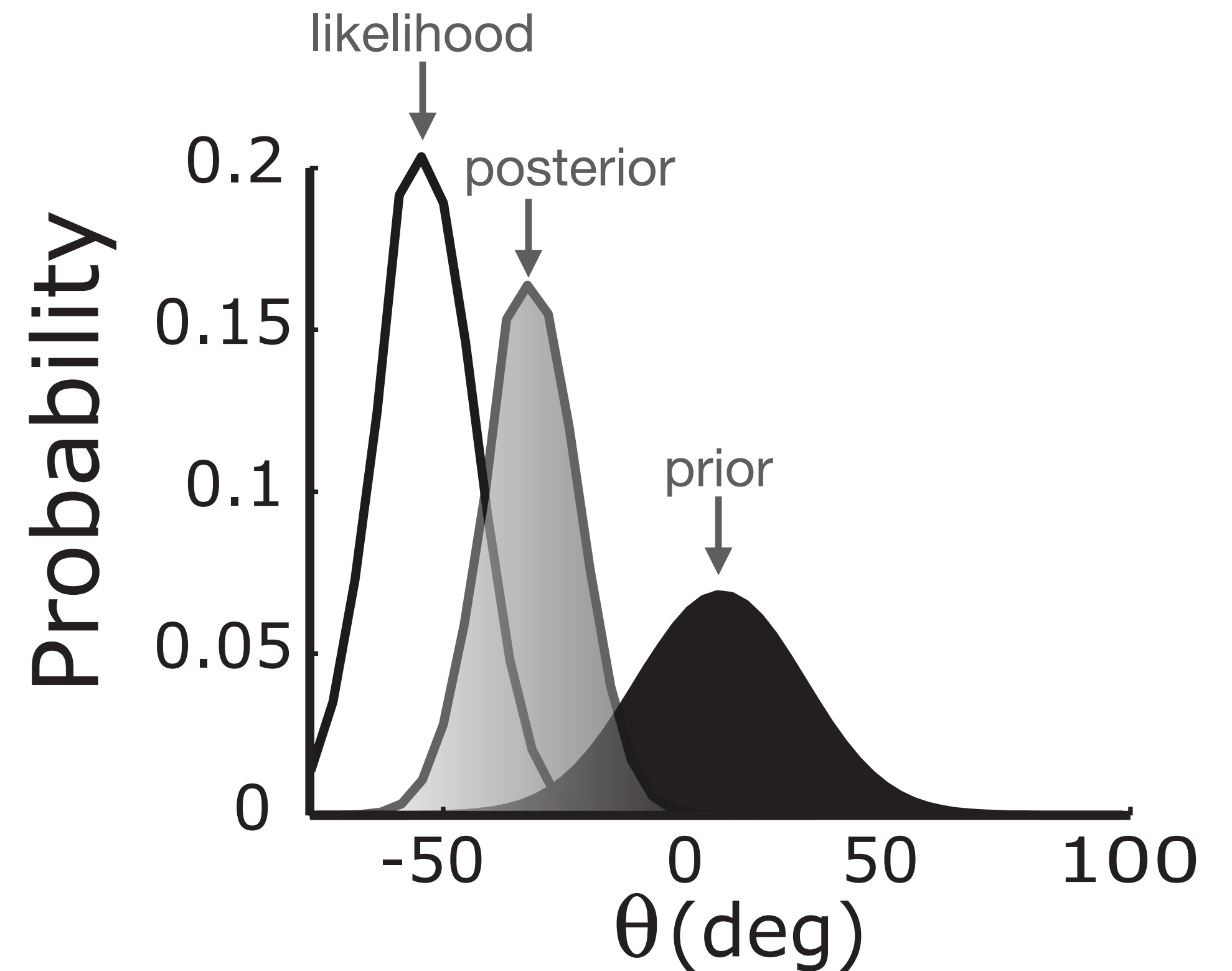
Discriminant analysis

Q: Which distribution does x_i come from?

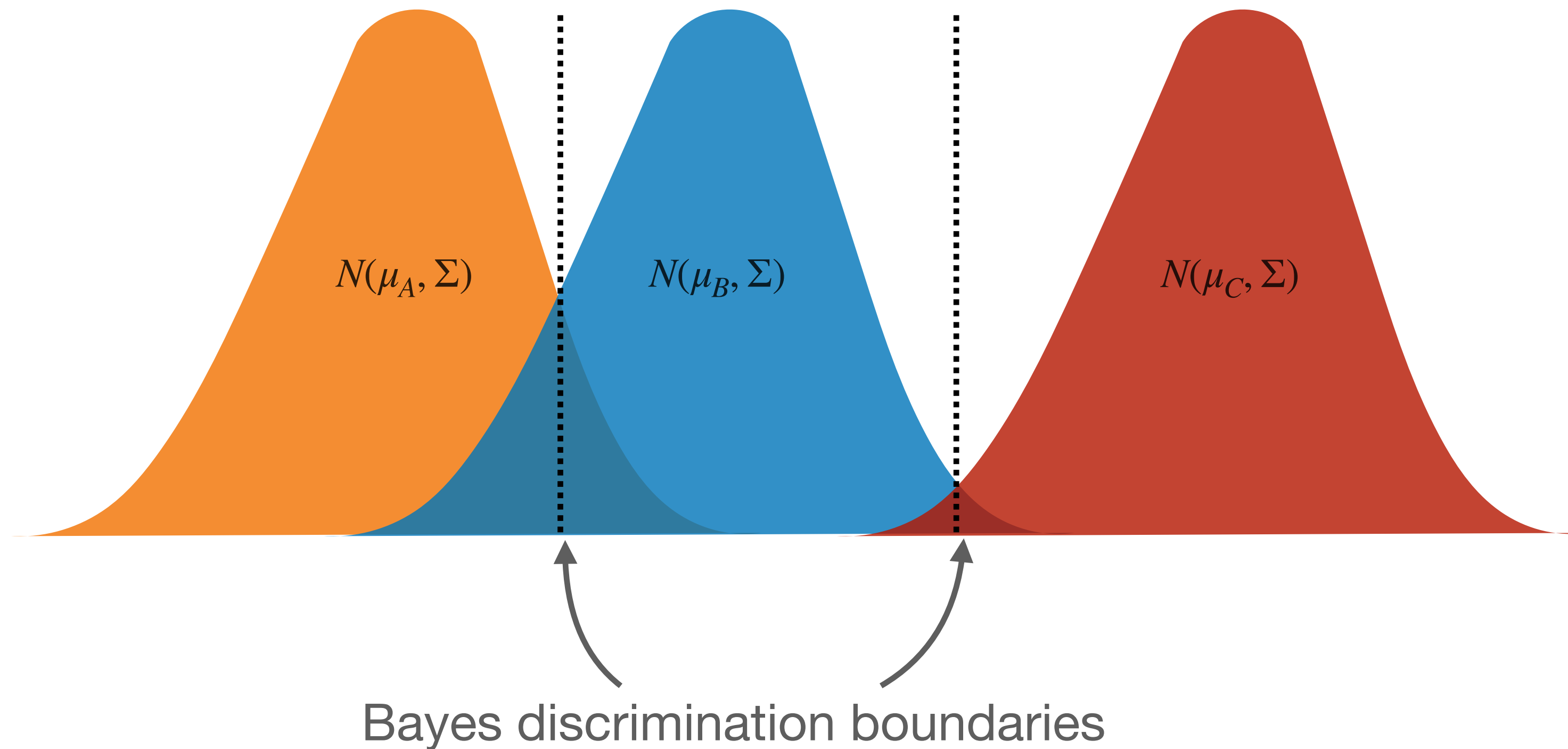


Bayes Theorem

$$\underbrace{P(A | B)}_{\text{posterior}} = \frac{\underbrace{P(B | A)}_{\text{likelihood}} \underbrace{P(A)}_{\text{prior}}}{\underbrace{P(B)}_{\text{marginal}}}$$



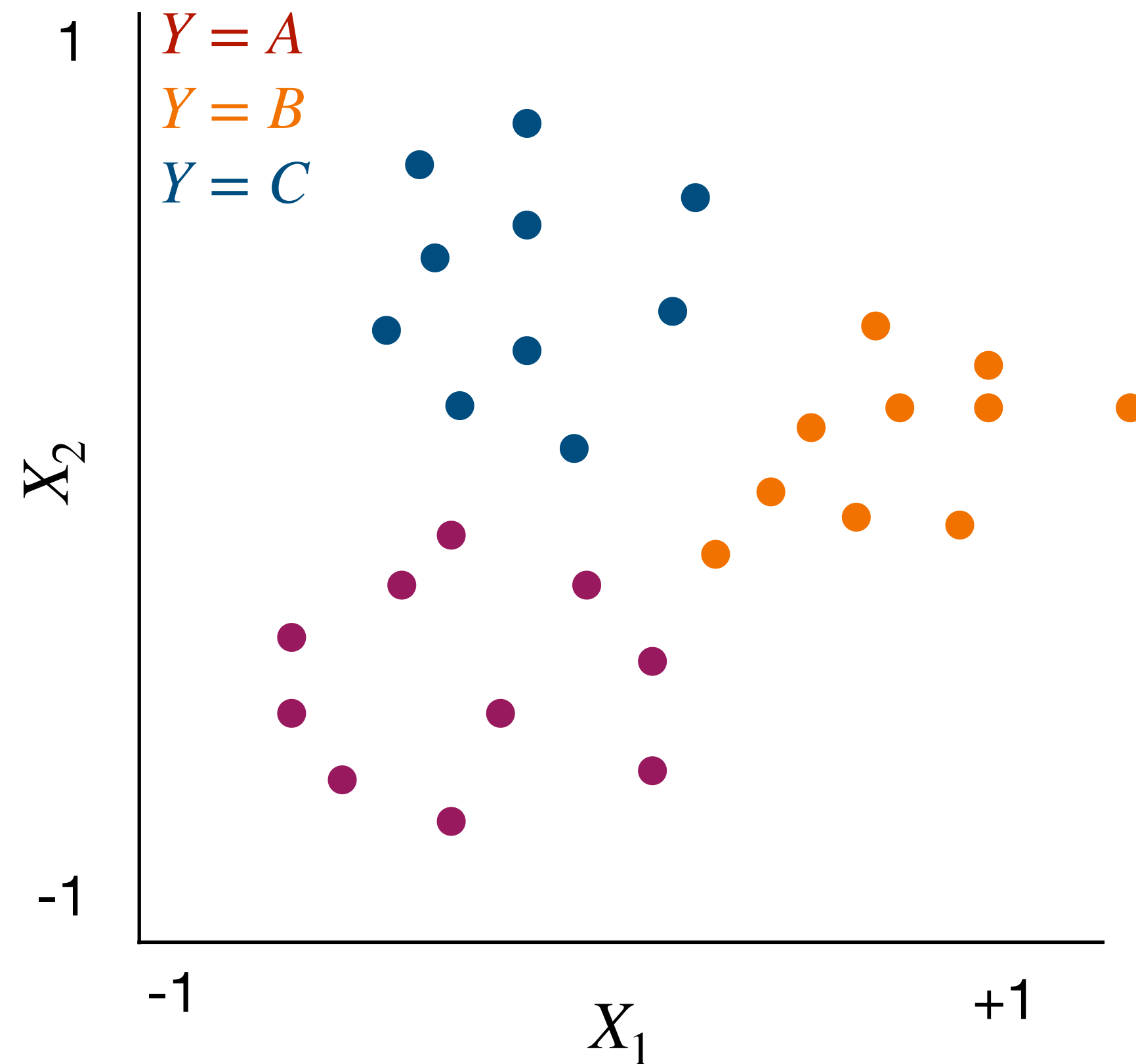
Bayesian classifier



$$P(Y = k | x_i) = \frac{P(x_i | Y = k)P(Y = k)}{P(x_i)}$$
$$\propto P(x_i | Y = k)P(Y = k)$$

Use probability theory to infer the separation of the underlying generative distributions.

Linear Discriminant Analysis (LDA)



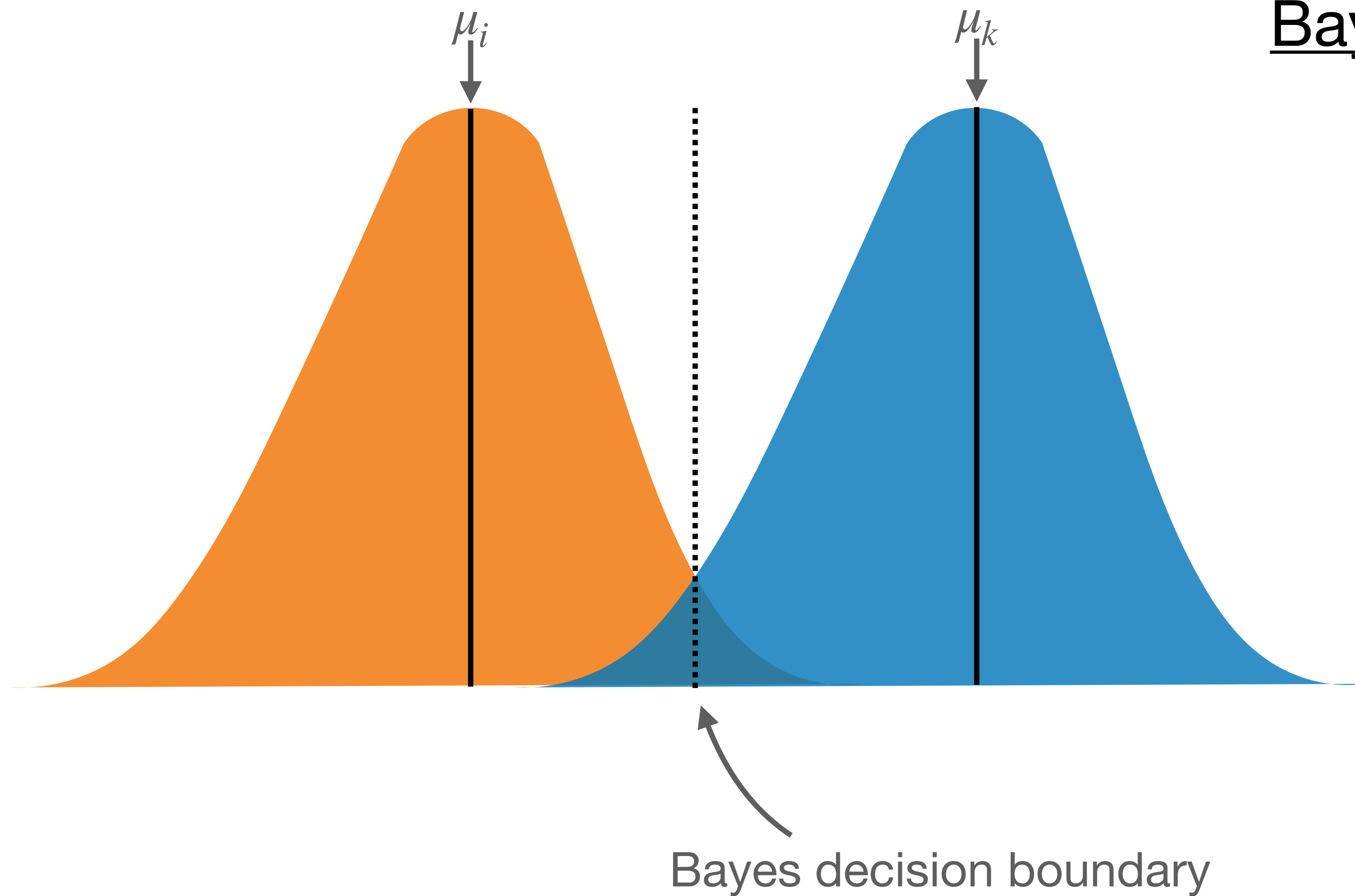
LDA

$\pi_i \sim$ prior for distribution i

$$P(Y = k | X = x_i) = P_k(x_i) = \frac{\pi_k f_k(x_i)}{\sum_{l=1}^p \pi_l f_l(x_i)}$$

$$P_k(x_i) = \frac{\pi_k \frac{1}{\Sigma \sqrt{2\pi}} e^{-\frac{(x_i - \mu_k)^2}{2\Sigma^2}}}{\sum_{l=1}^p \pi_l \frac{1}{\Sigma \sqrt{2\pi}} e^{-\frac{(x_i - \mu_l)^2}{2\Sigma^2}}}$$

Optimal separation of Gaussians

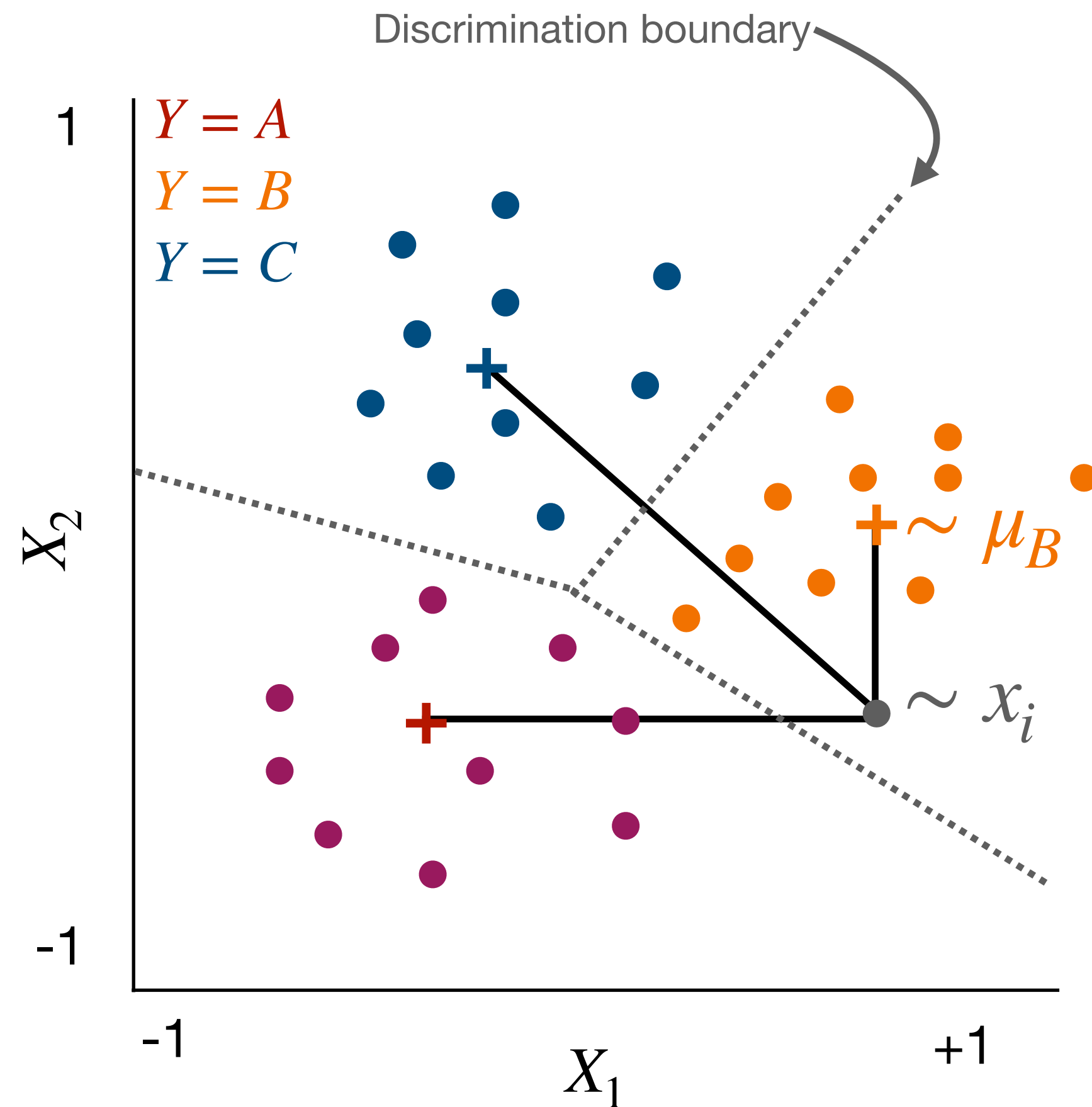


Bayes Decision Boundary

$$\frac{\mu_i^2 - \mu_k^2}{2(\mu_i - \mu_k)} = \frac{\mu_i + \mu_k}{2}$$

Observations on one side of the boundary are group i , observations on the other are k .

Group discriminations



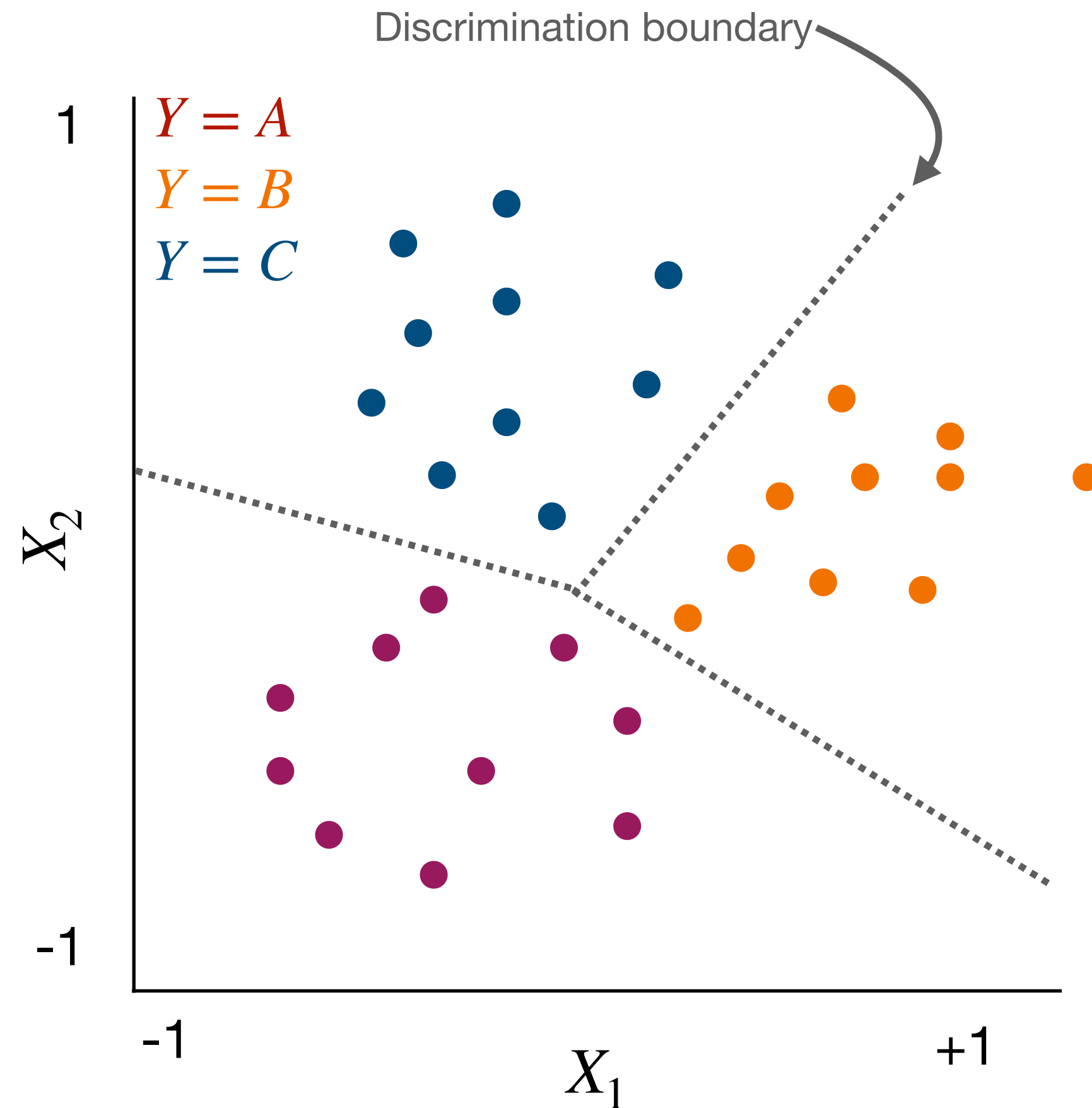
Discriminant function

$$\delta_k(x_i) = x_i \frac{\mu_k}{\Sigma^2} - \frac{\mu_k^2}{2\Sigma^2} - \ln(\pi_k)$$

$\pi_k = \frac{n_k}{n_{all}}$

$\uparrow \delta_k$ means x_i more likely arises from distribution k .

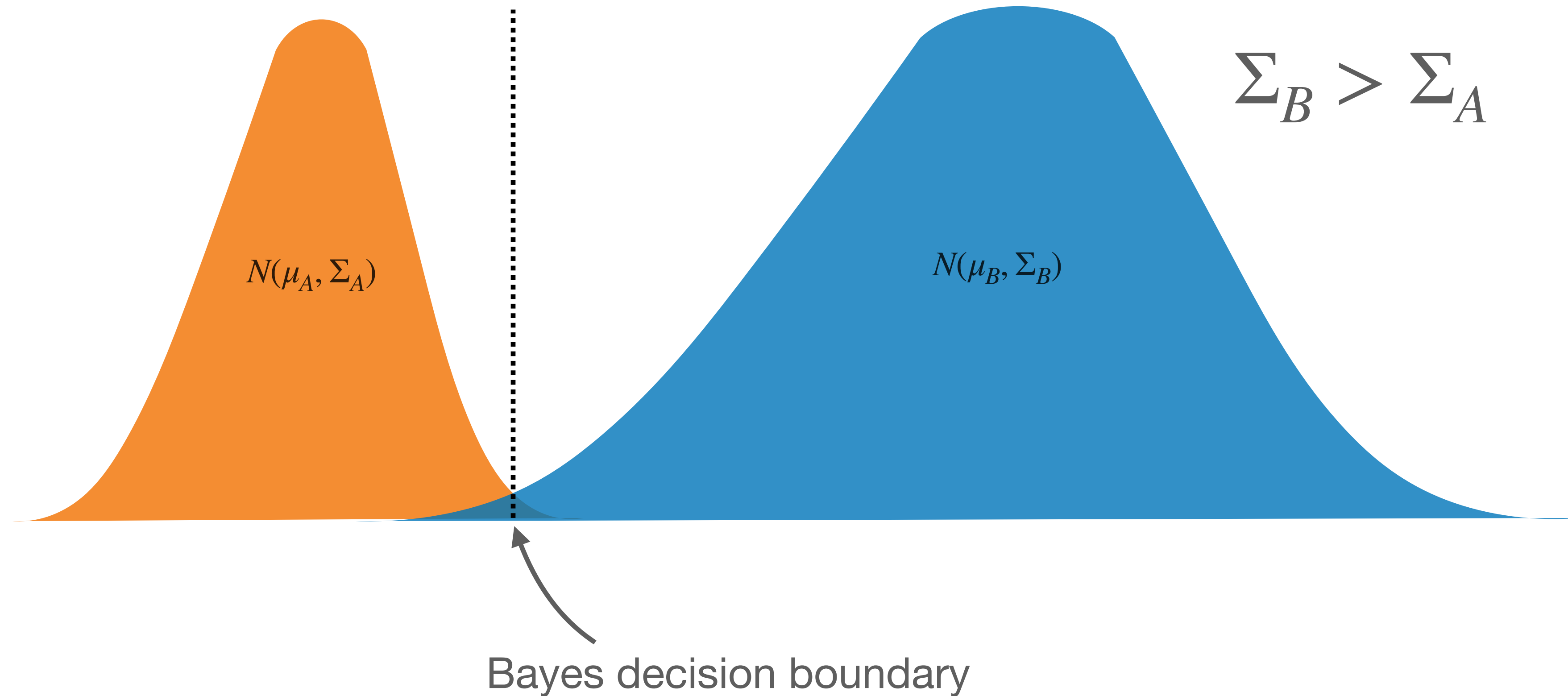
Linear Discriminant Analysis (LDA)



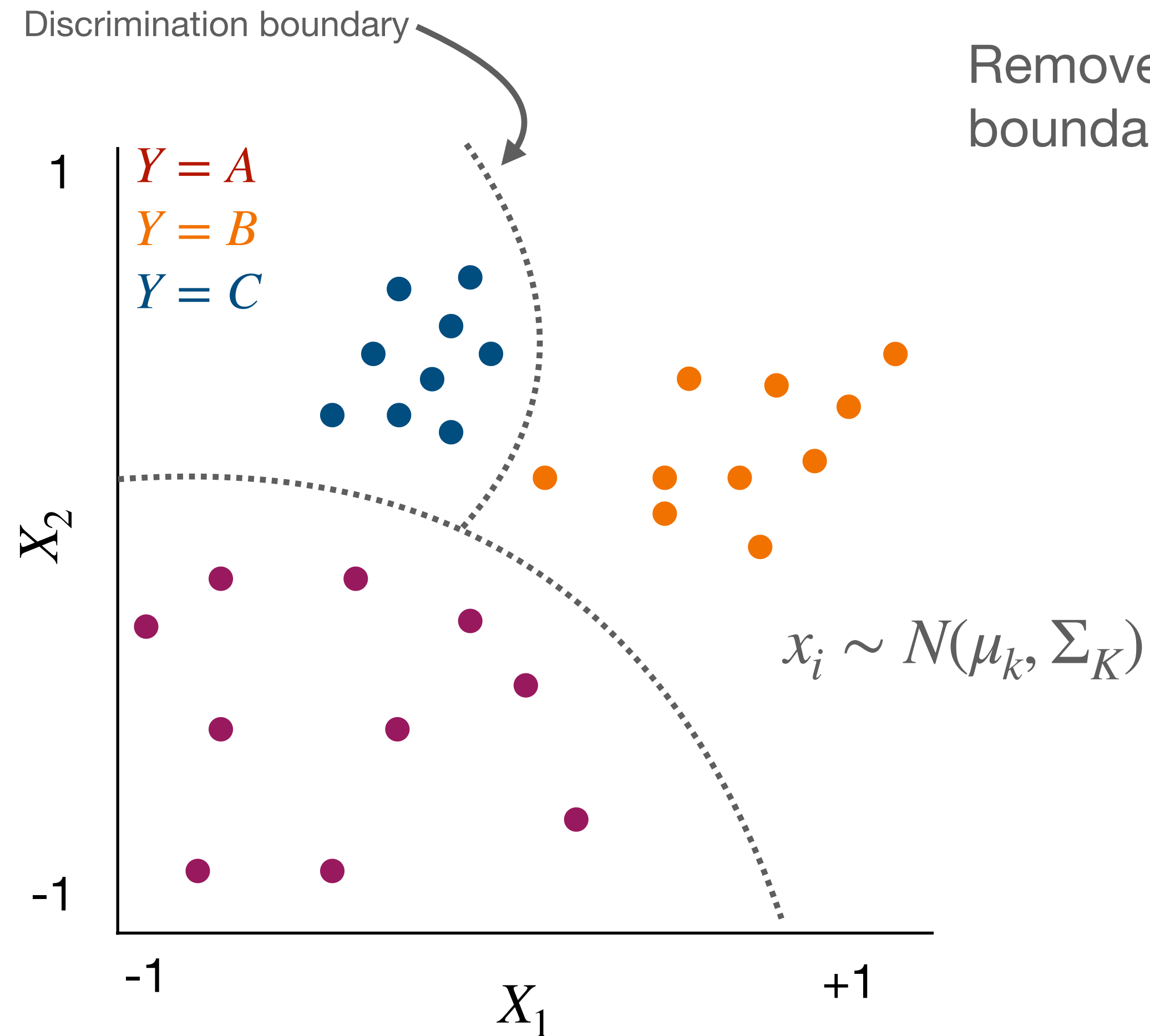
Assumptions:

1. x_i is a multivariate normal random variable.
2. All generative distributions have the same variance Σ^2 .
3. There is no collinearity between groups.
4. Errors are independent.

What if variances are not the same?



Quadratic Discriminant Analysis (QDA)

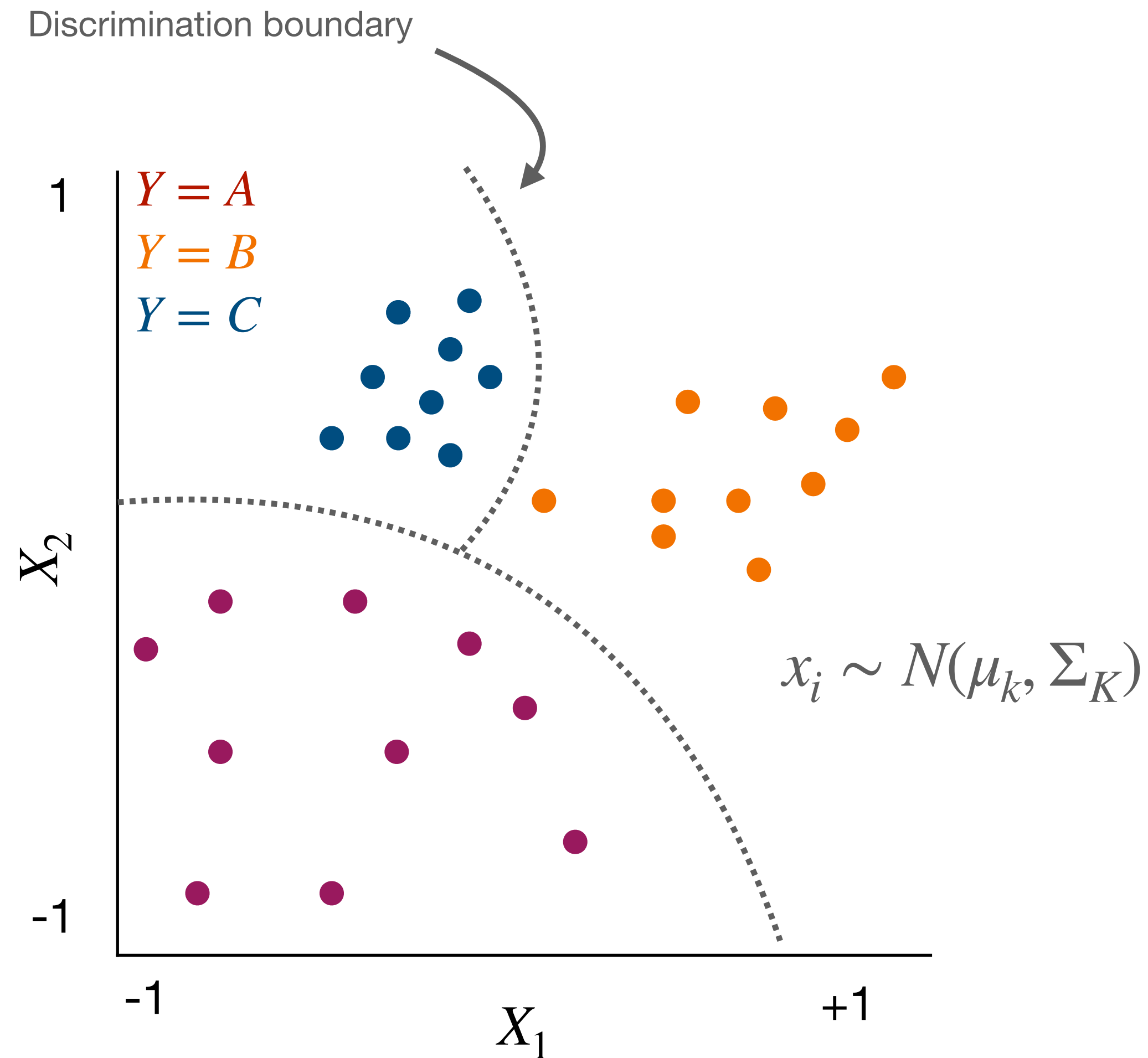


Remove the assumption of equal variances and the decision boundaries become curves, instead of straight lines.

Discriminant function

$$\delta_k(x_i) = -\frac{(x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k)}{2} - \frac{1}{2} \ln(\Sigma_k) + \ln(\pi_k)$$

Quadratic Discriminant Analysis (QDA)



Assumptions:

1. x_i is a multivariate normal random variable.
2. There is no collinearity between groups.
3. Errors are independent.

Take home message

- While logistic regression offers an intuitive extension to linear regression for classification problems, discriminant analyses are more flexible for classification in cases where your data consists of more than 2 categories.