

# Generalized Linear Models (GLMs)

Dr. Alessandro Filazzola & Dr. Christopher Lortie

BIOL 5018

# Outline & Learning Outcomes

## **Part A – Conducting GLMs**

- Describe the uses of GLMs relative to linear regression
- Be able to conduct a GLM in R
- Gain a basic understanding of the mechanics and outputs of a GLM

## **Part B – Prediction**

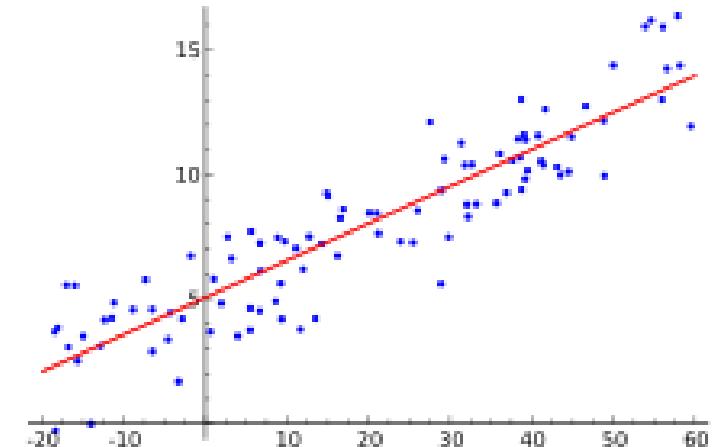
- Describe how GLMs can be used for predictions with new data
- Explain the approach for optimizing predictions

## **Part C – Post-hoc analyses and mixed models**

- Understand the differences between GLMs and GLMMs
- Explain the purpose of post-hoc analyses and when they may be used
- Be able to conduct a GLMM and post-hoc analysis in R

# Generalized Linear Models (GLM)

- General linear models = typical linear regression
- Generalized Linear Models = distribution + link function
  - Poisson regression
  - Negative binomial regression
  - Gamma regression
  - Zero-inflated regression
  - Logistic regression (sometimes referred to as binomial)



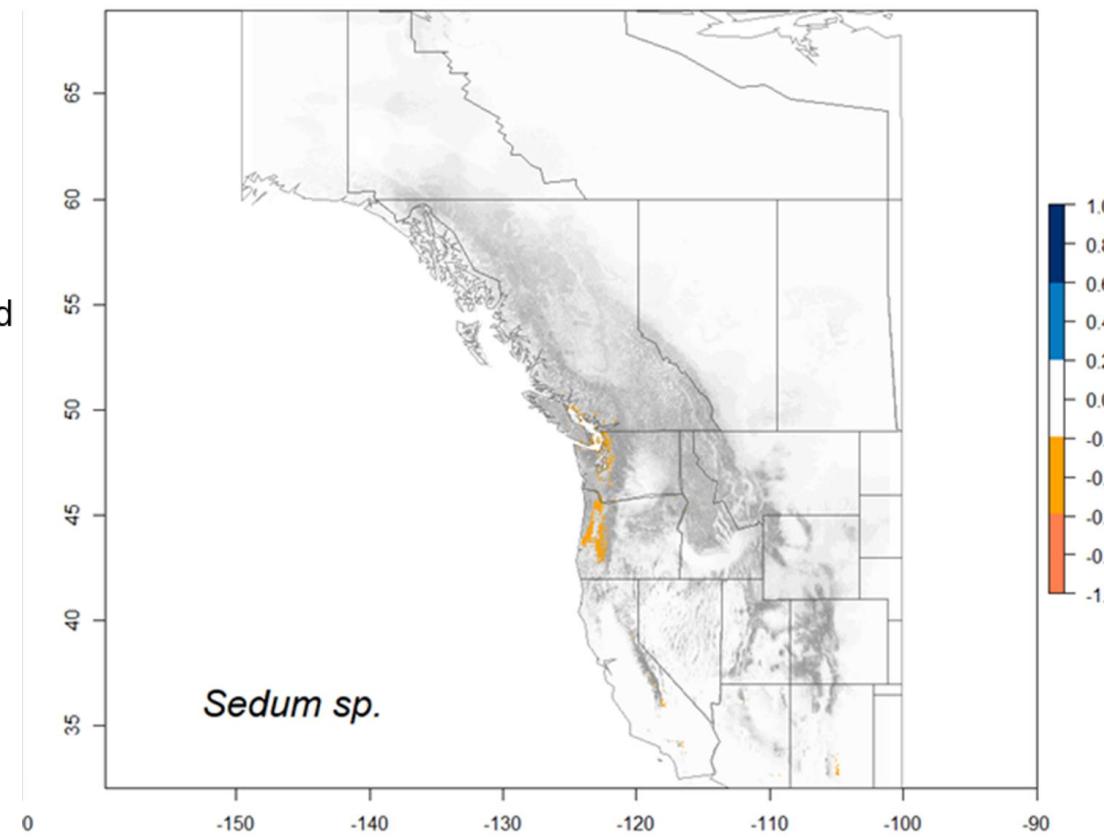
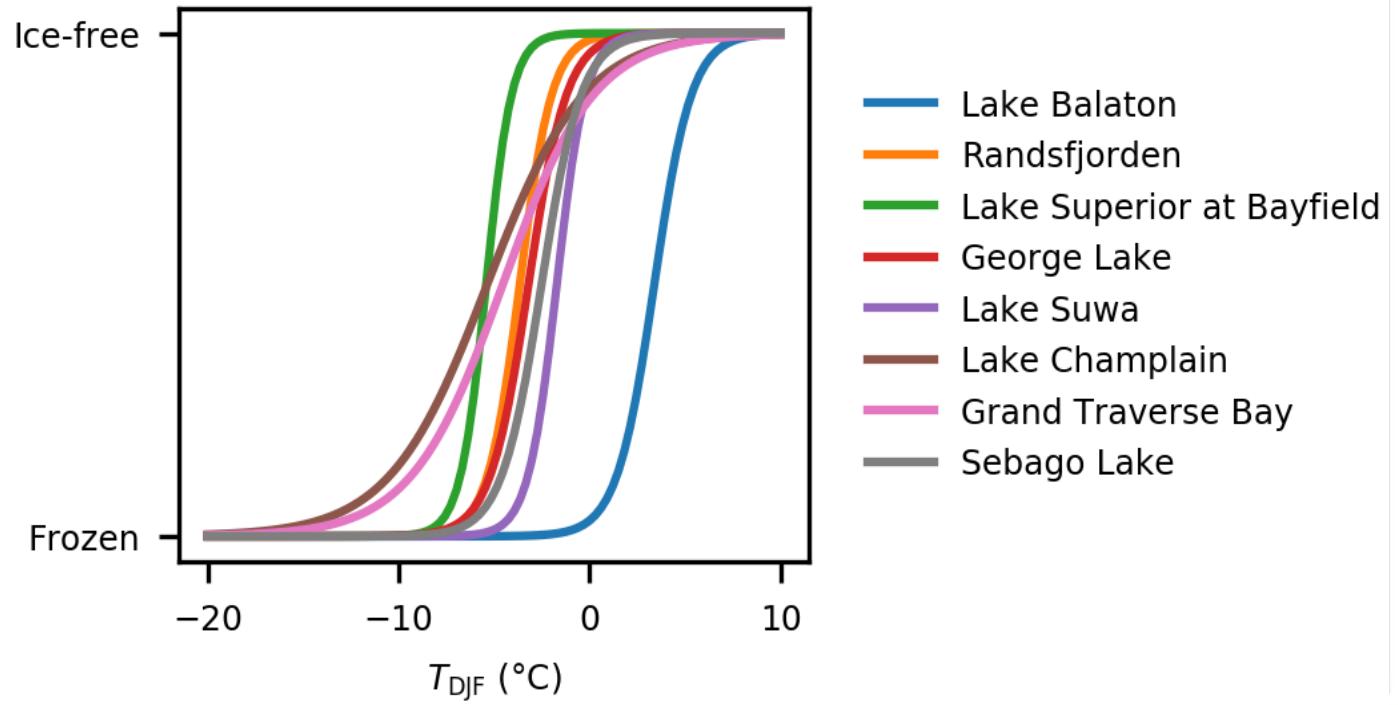
# Why use a GLM?

- 1) The response variable is not normally distributed, or the range is restricted
- 2) The variance of the response variable depends on the mean
- 3) GLMs can fit to particular distributions without transformations.  
The benefits of this include i) the homogeneity of variance does NOT need to be satisfied and ii) the response variable does not need to be changed.

$$E(\log(Y)) \neq \log(E(Y))$$

# Examples of GLMs

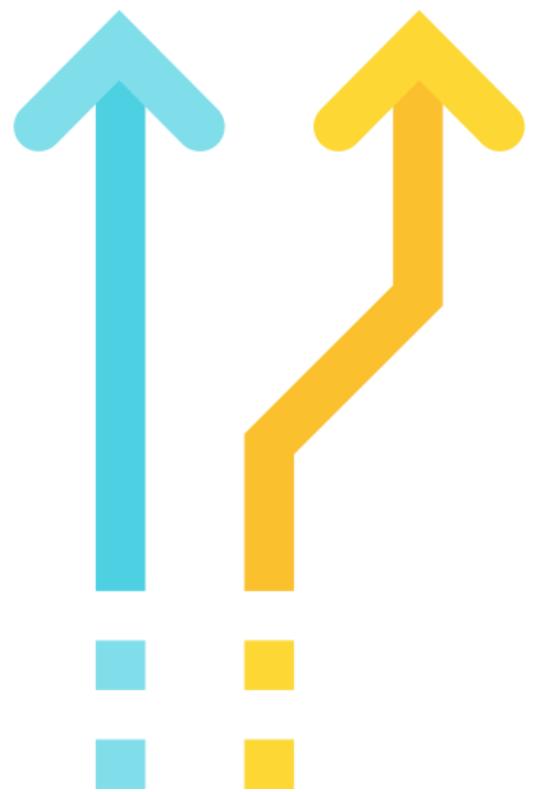
Used frequently to **test** and **predict**



# Parallels with linear regression

Both logistic and linear regression can:

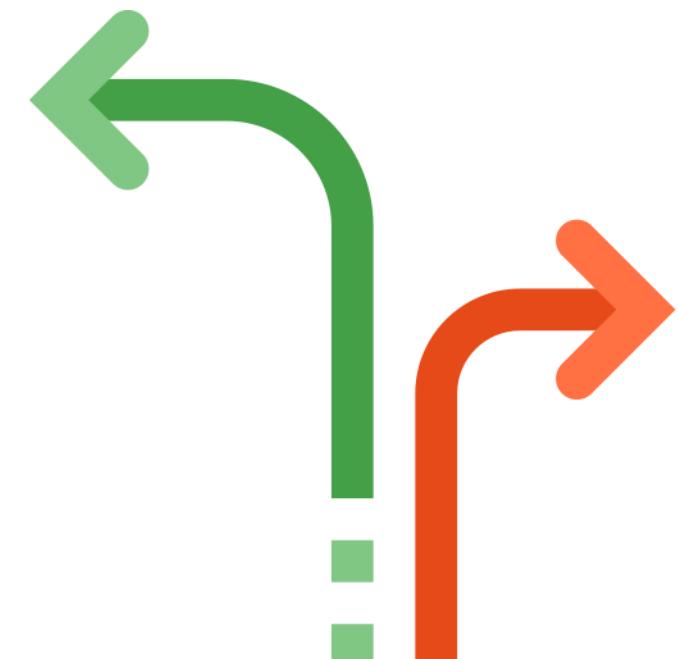
- Test for significance of a predictor
- Calculate an effect size
- Use continuous or categorical predictors
- Can predict values for new values



# Differences with linear regression

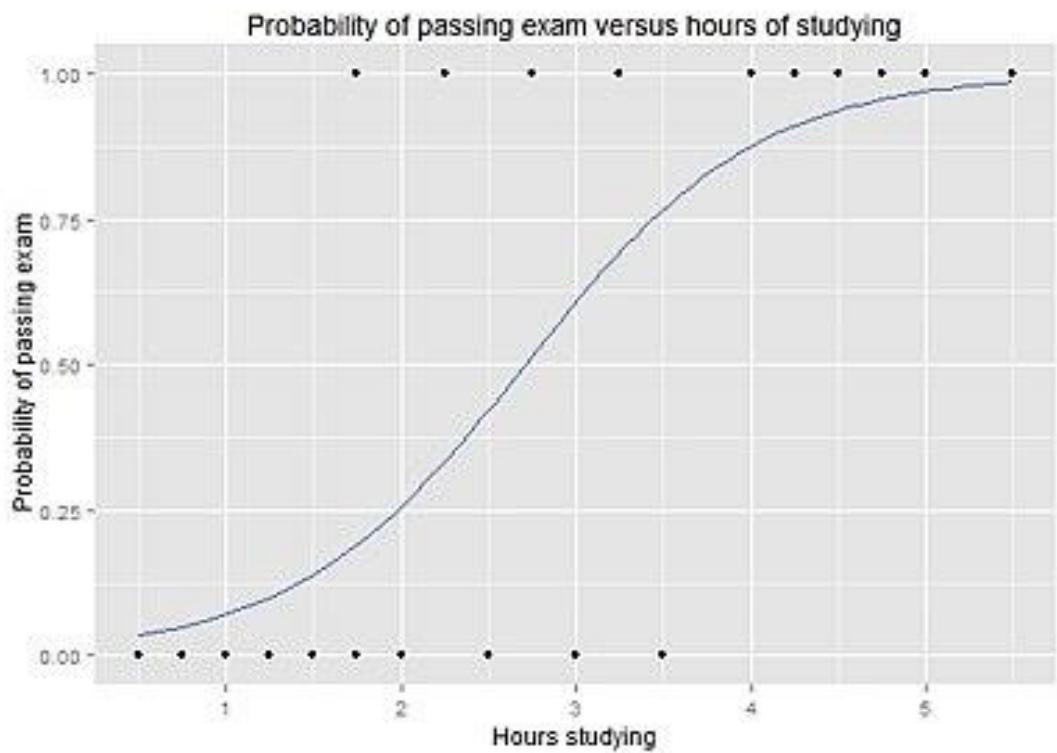
Logistic regression differs from linear regression in that:

- Response variable can be complex (e.g., binary, discrete, continuous)
- There is no “true” R<sup>2</sup> or residuals
- Uses Maximum Likelihood instead of Sum of Squares

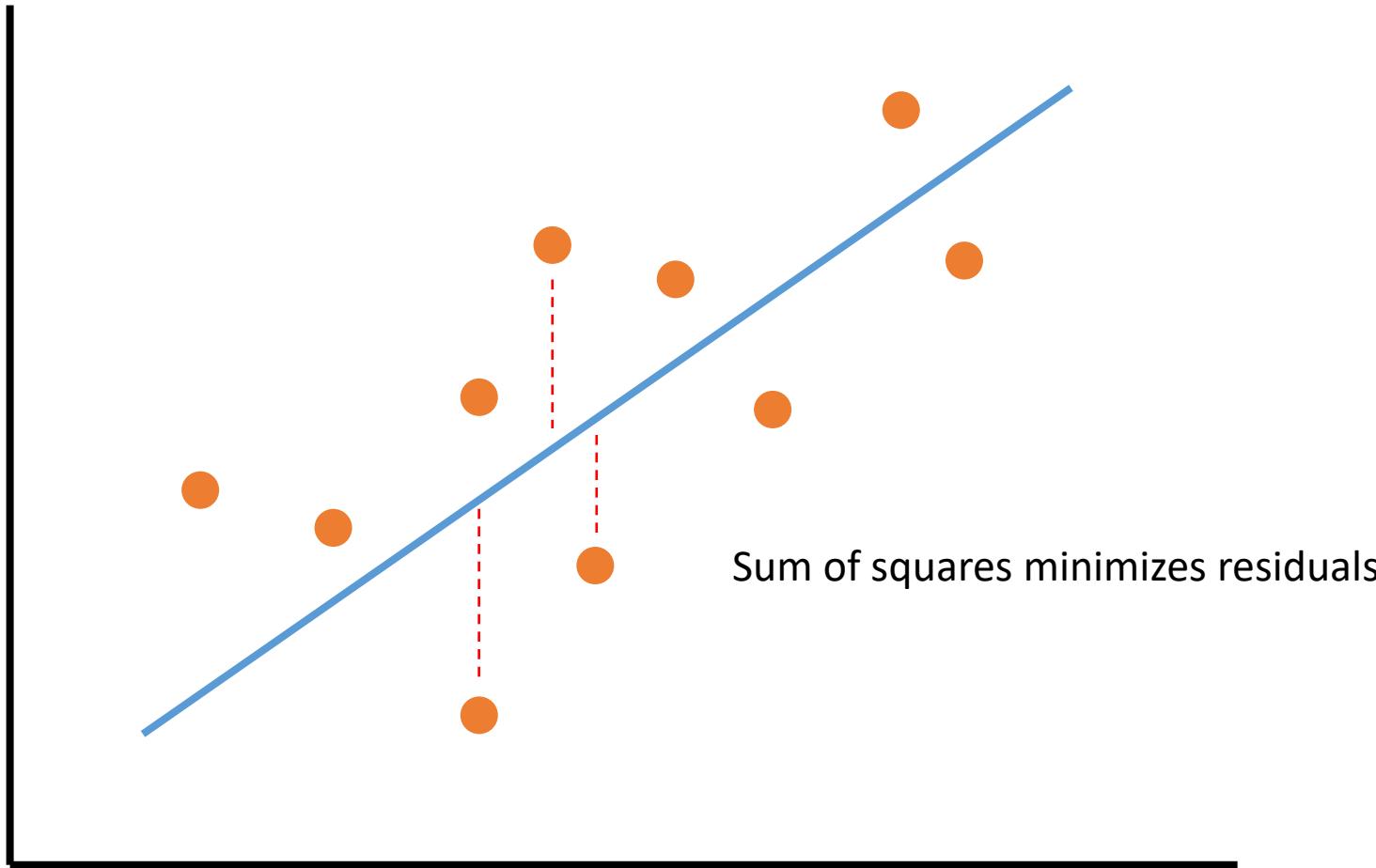


# Understanding GLMs

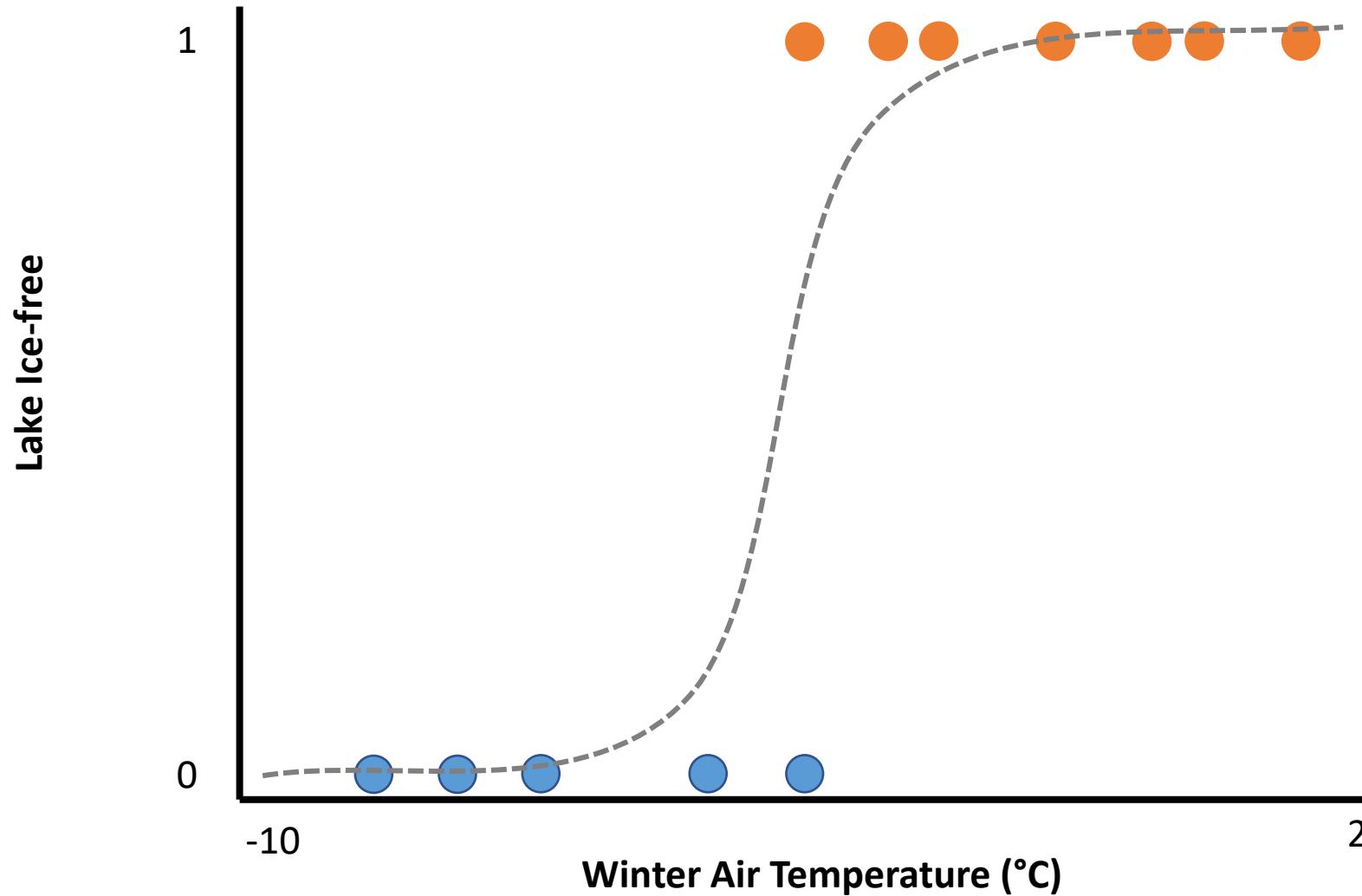
- 1) Fitting a line
- 2) The outputs
- 3) Calculating fit



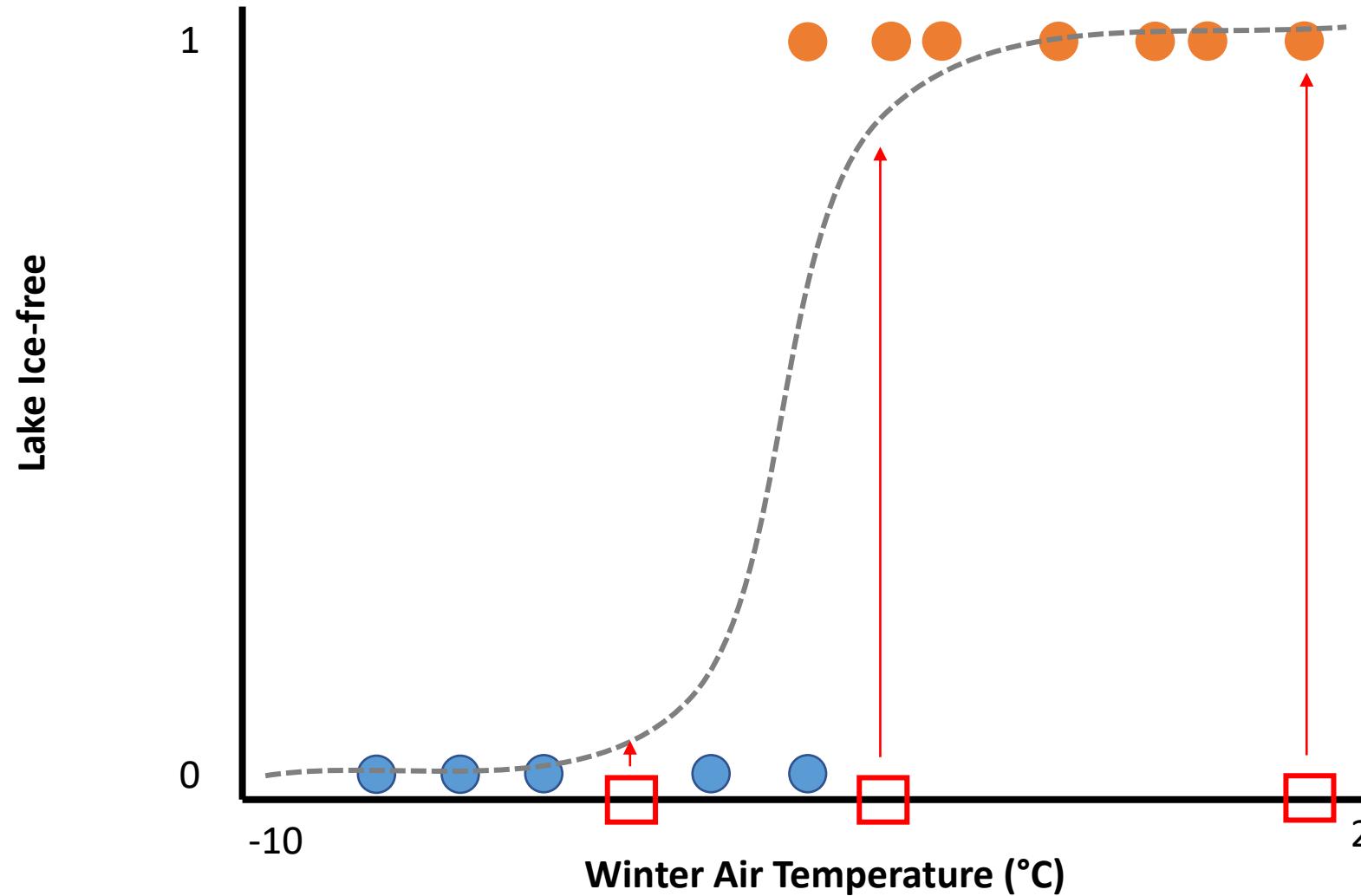
# 1) Fitting a line: linear regression



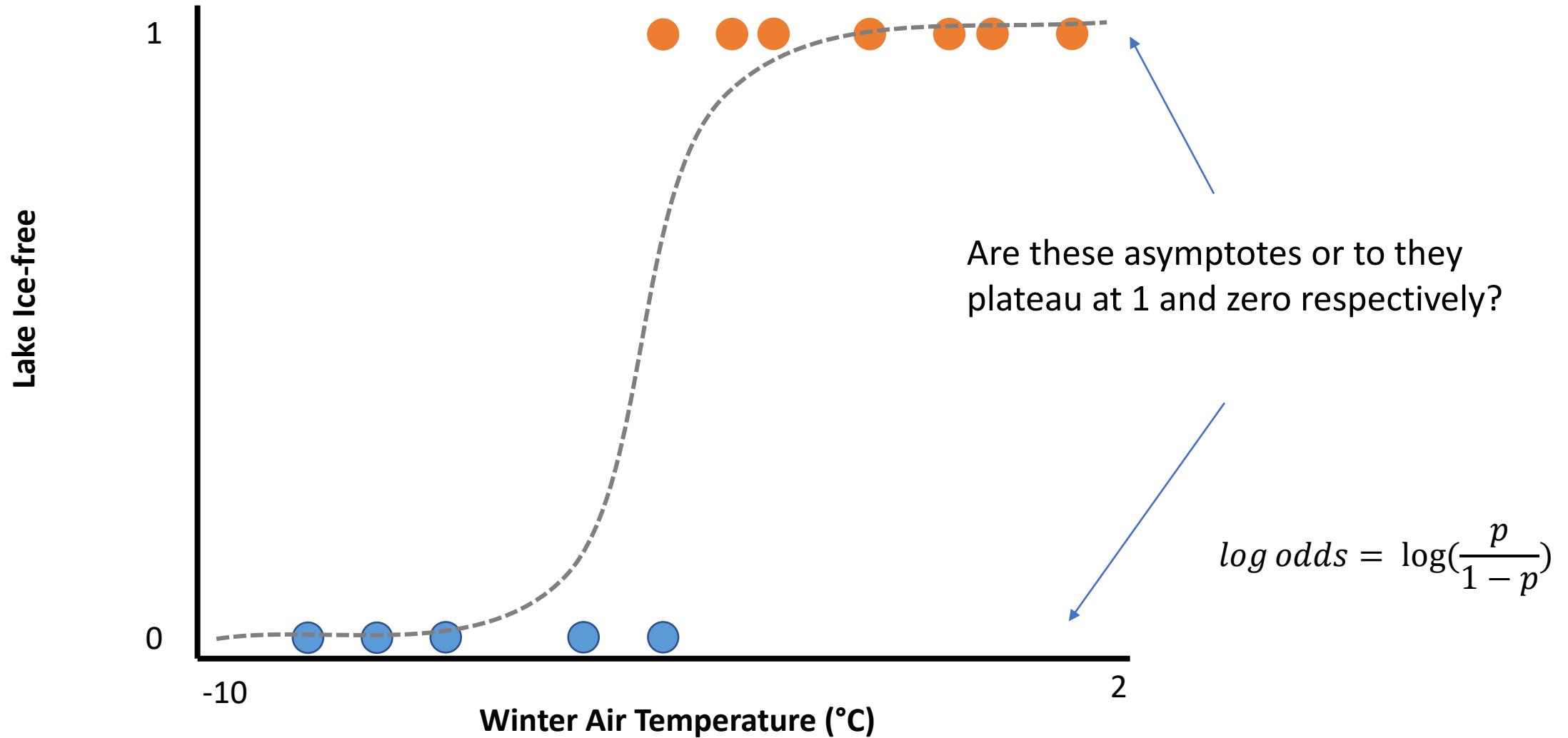
# 1) Fitting a line: logistic regression



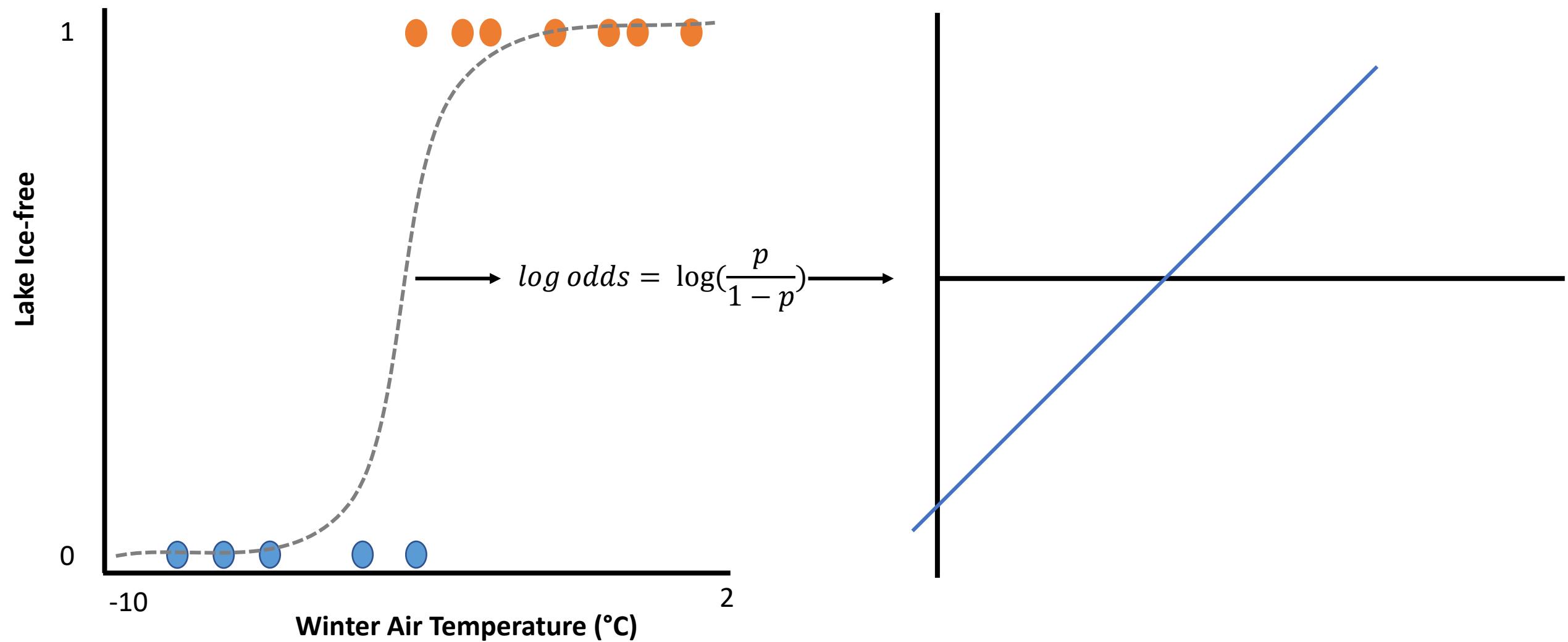
# 1) Fitting a line: logistic regression



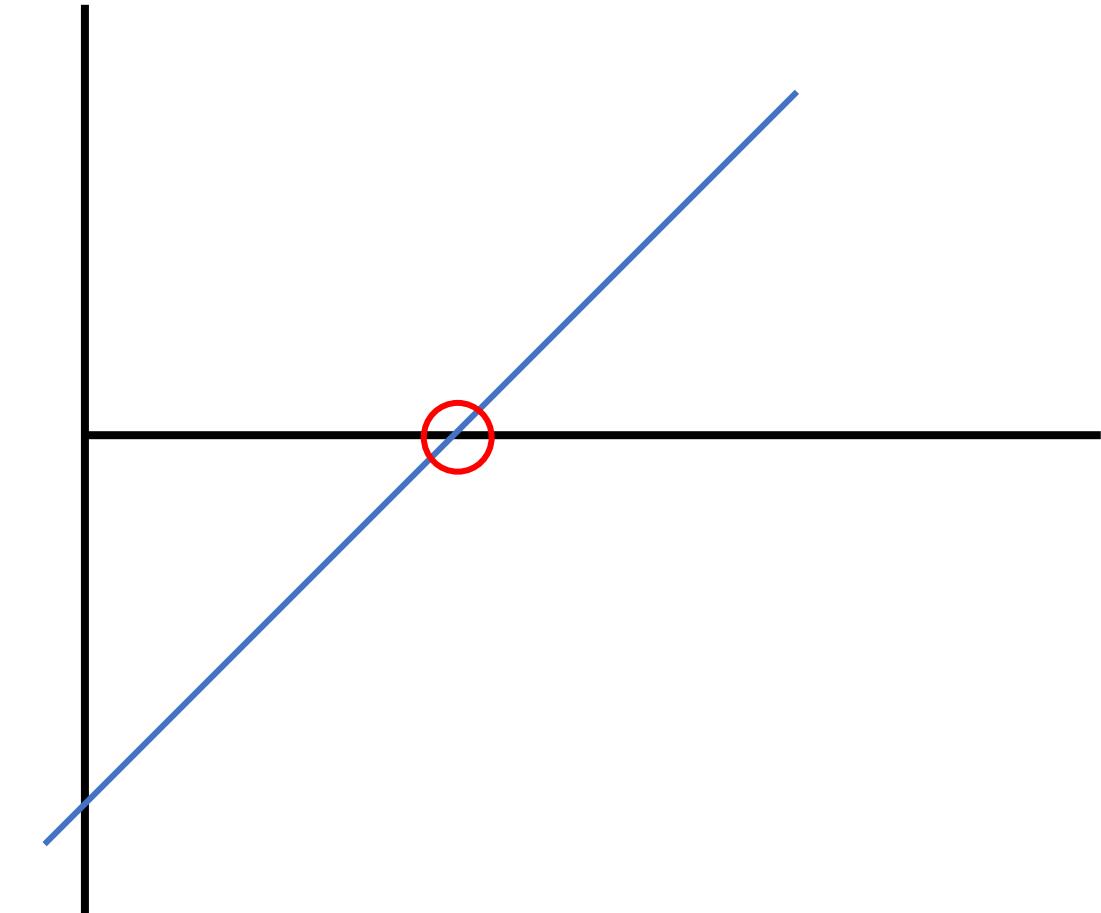
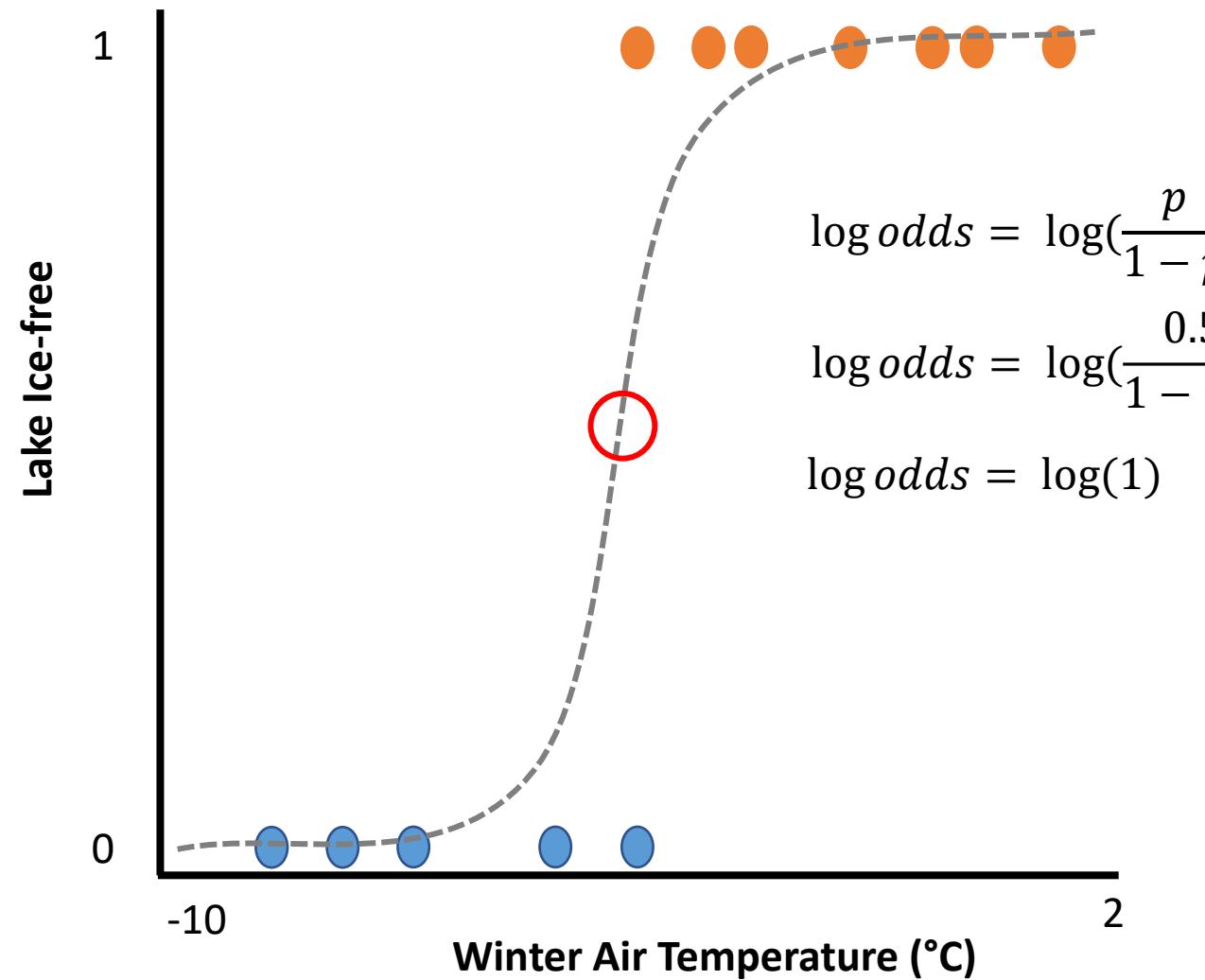
# Question



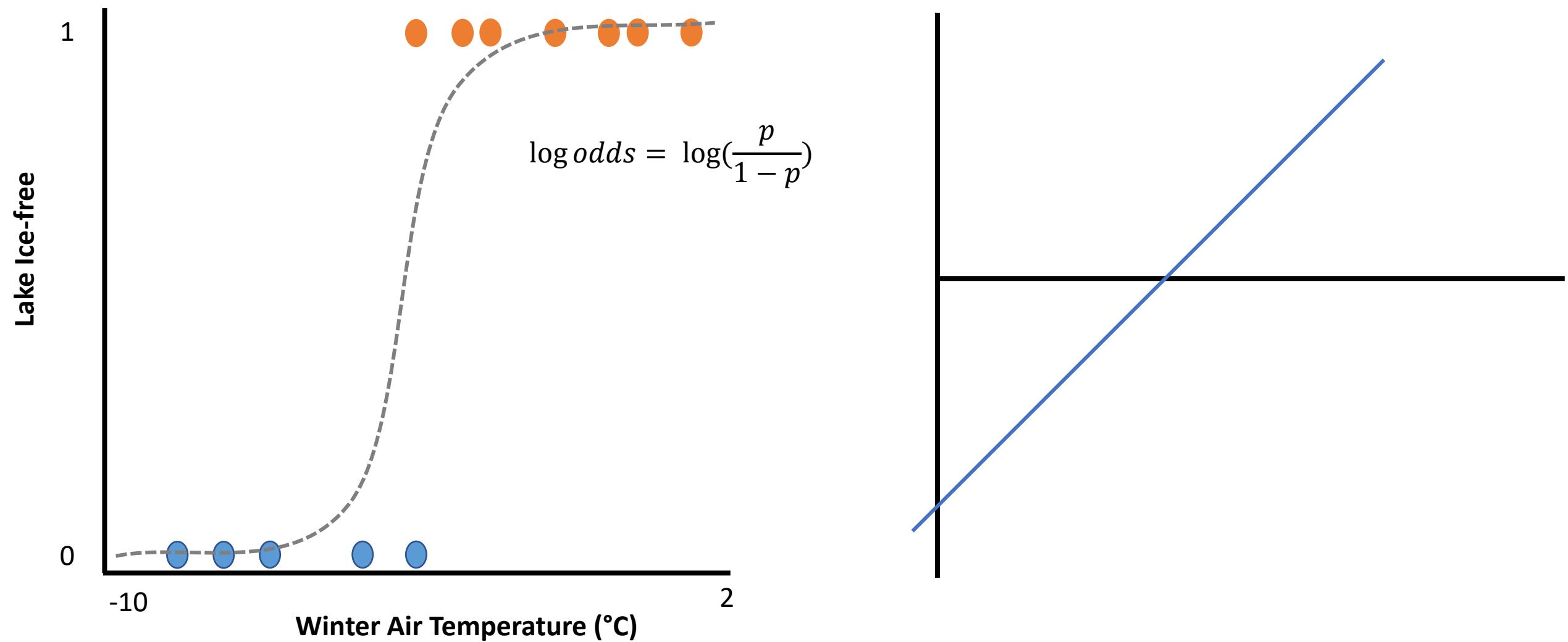
# 1) Putting the *linear* in generalized linear model



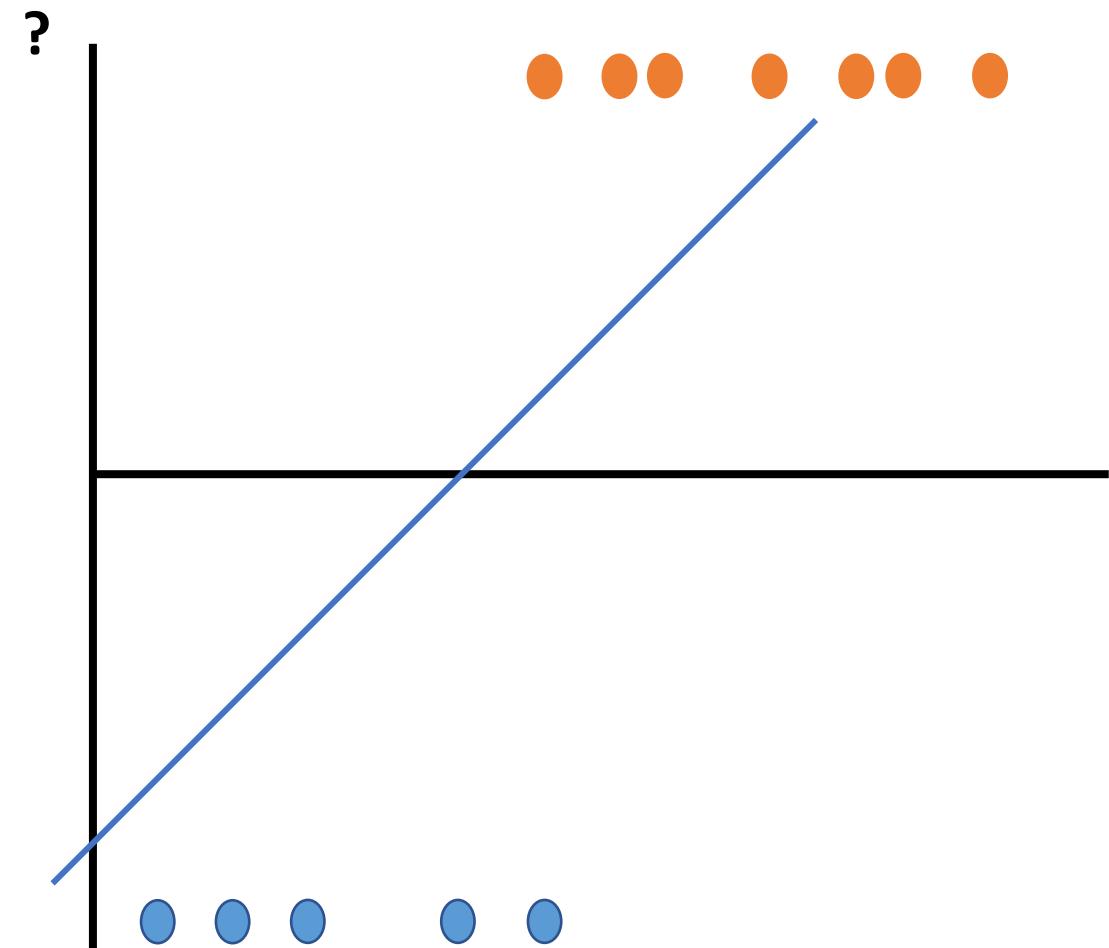
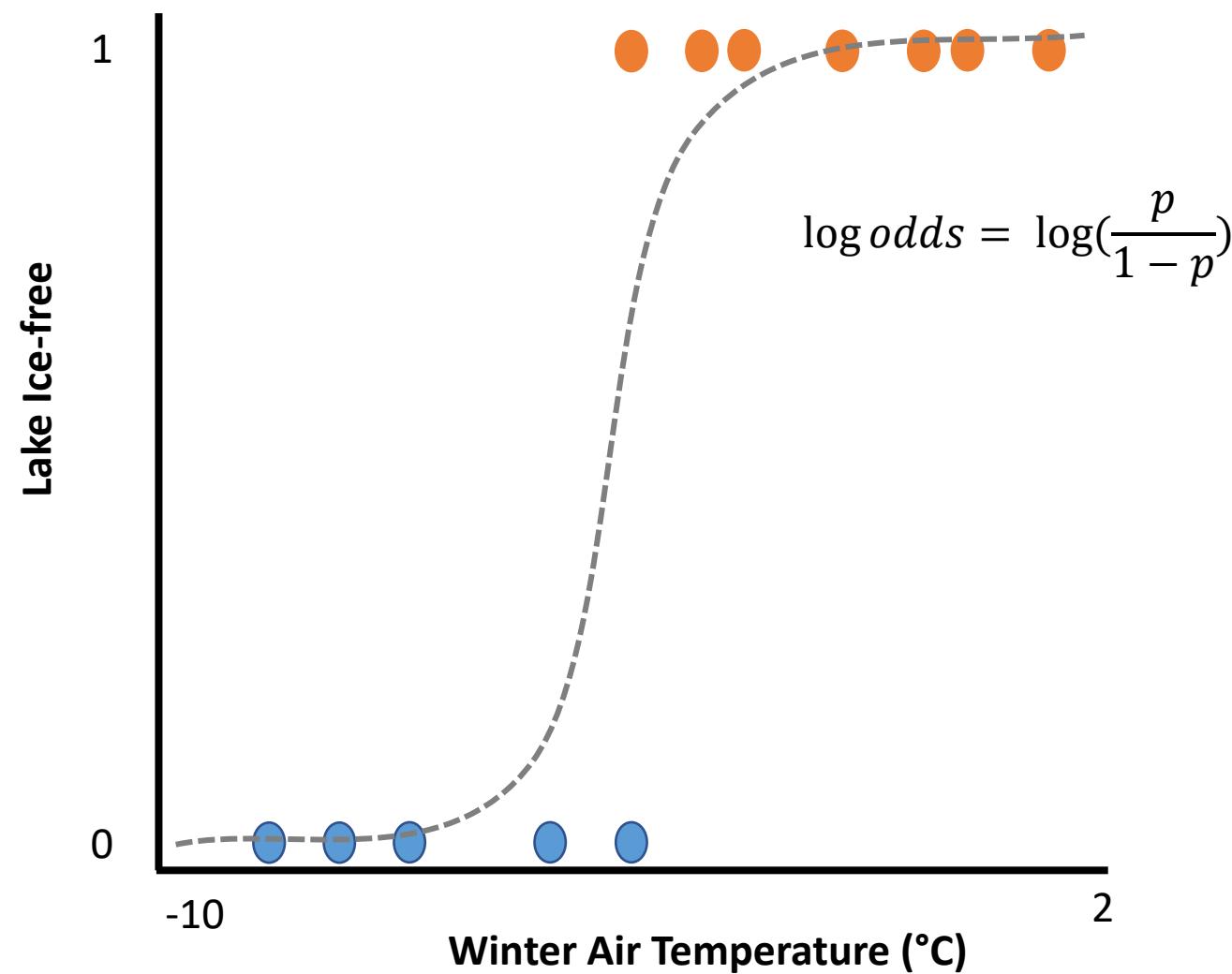
# 1) Putting the *linear* in generalized linear model



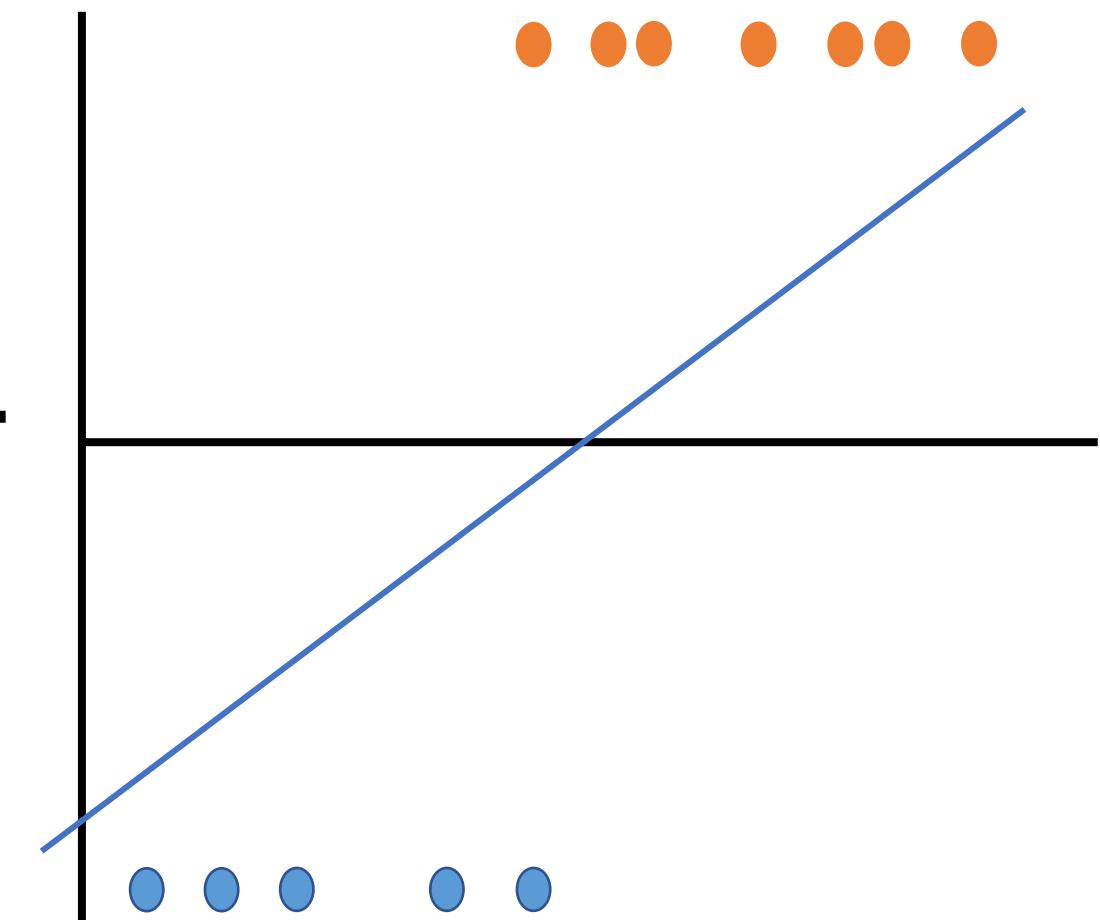
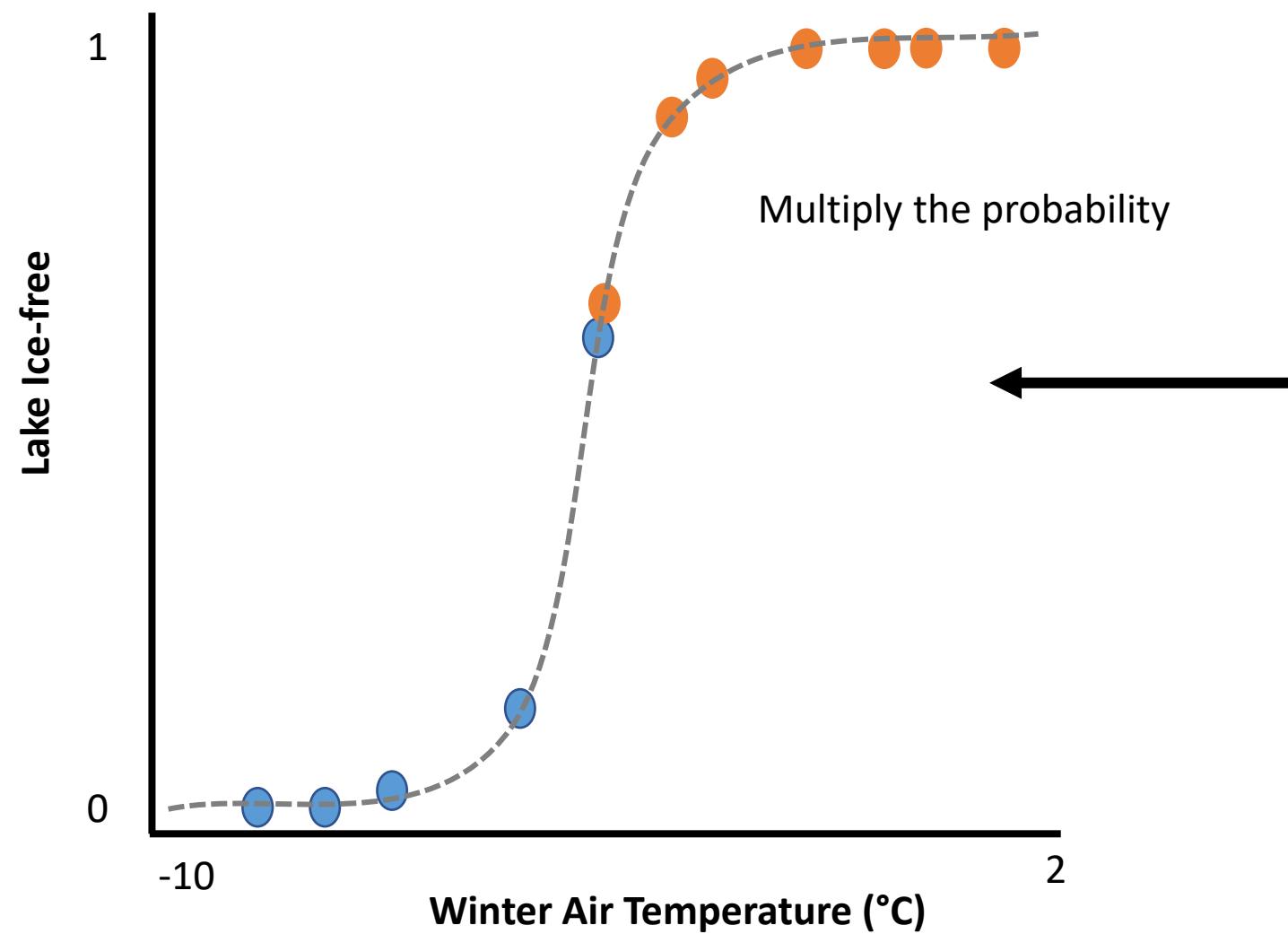
# 1) Putting the *linear* in generalized linear model



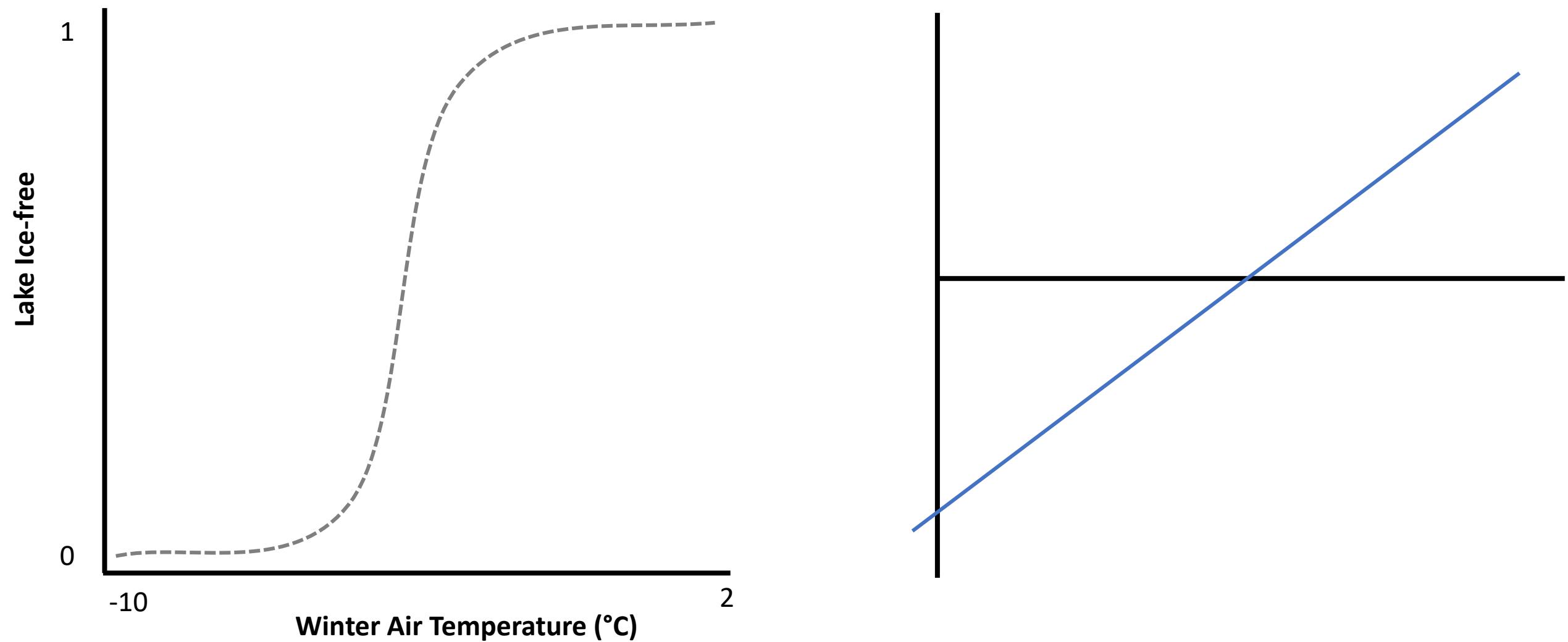
# 1) Putting the *linear* in generalized linear model



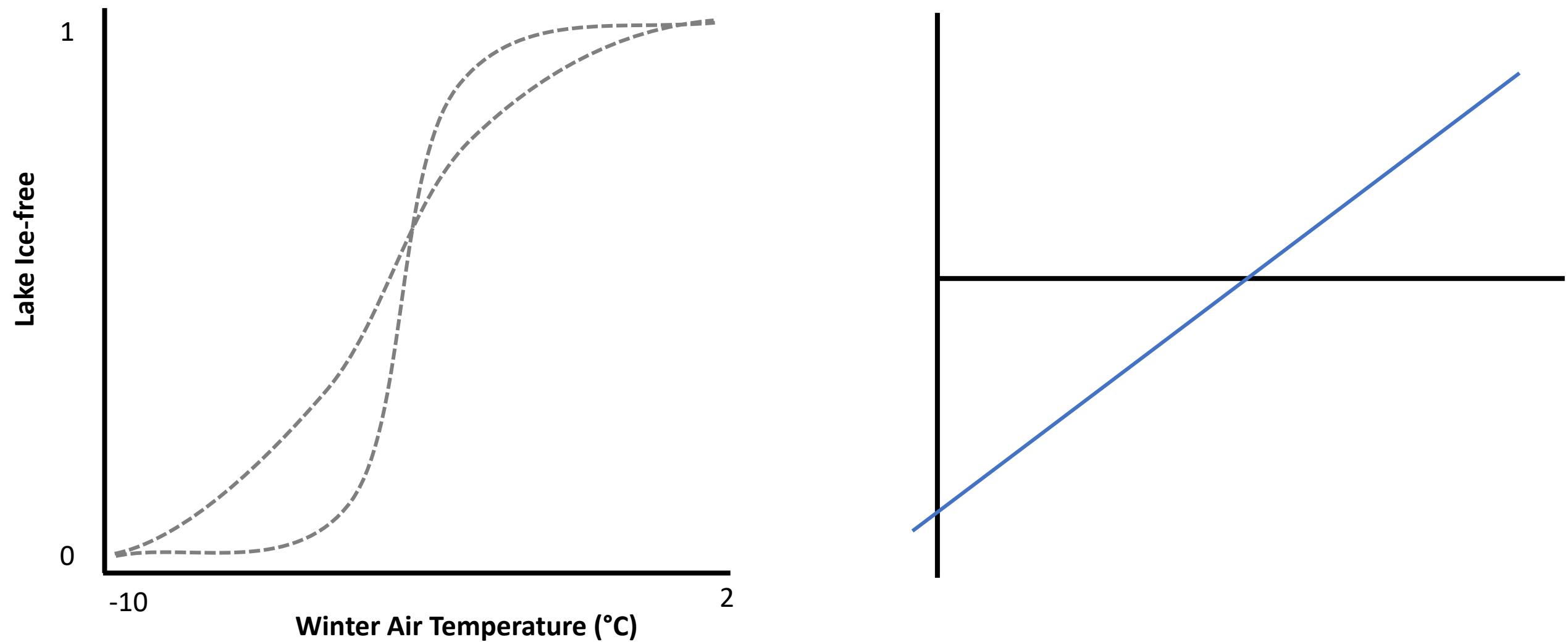
# 1) Fitting the model



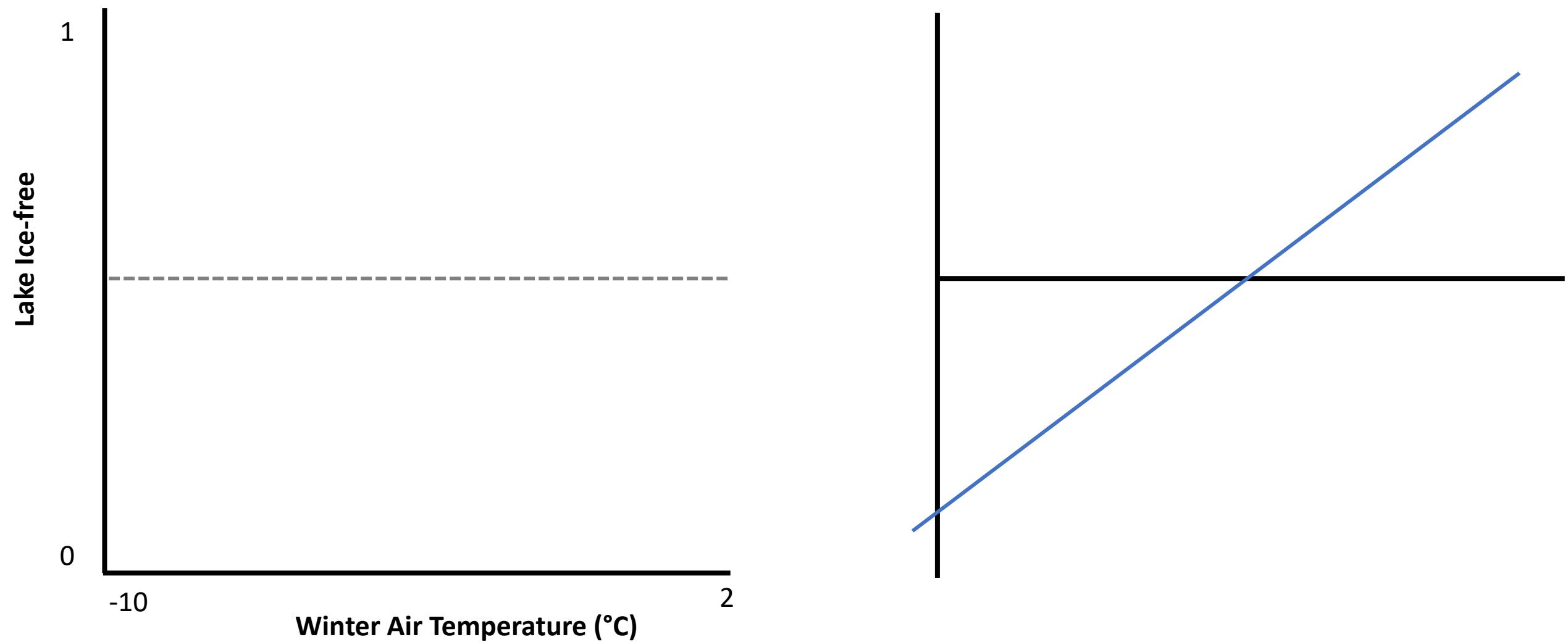
# 1) Fitting the model – maximum likelihood



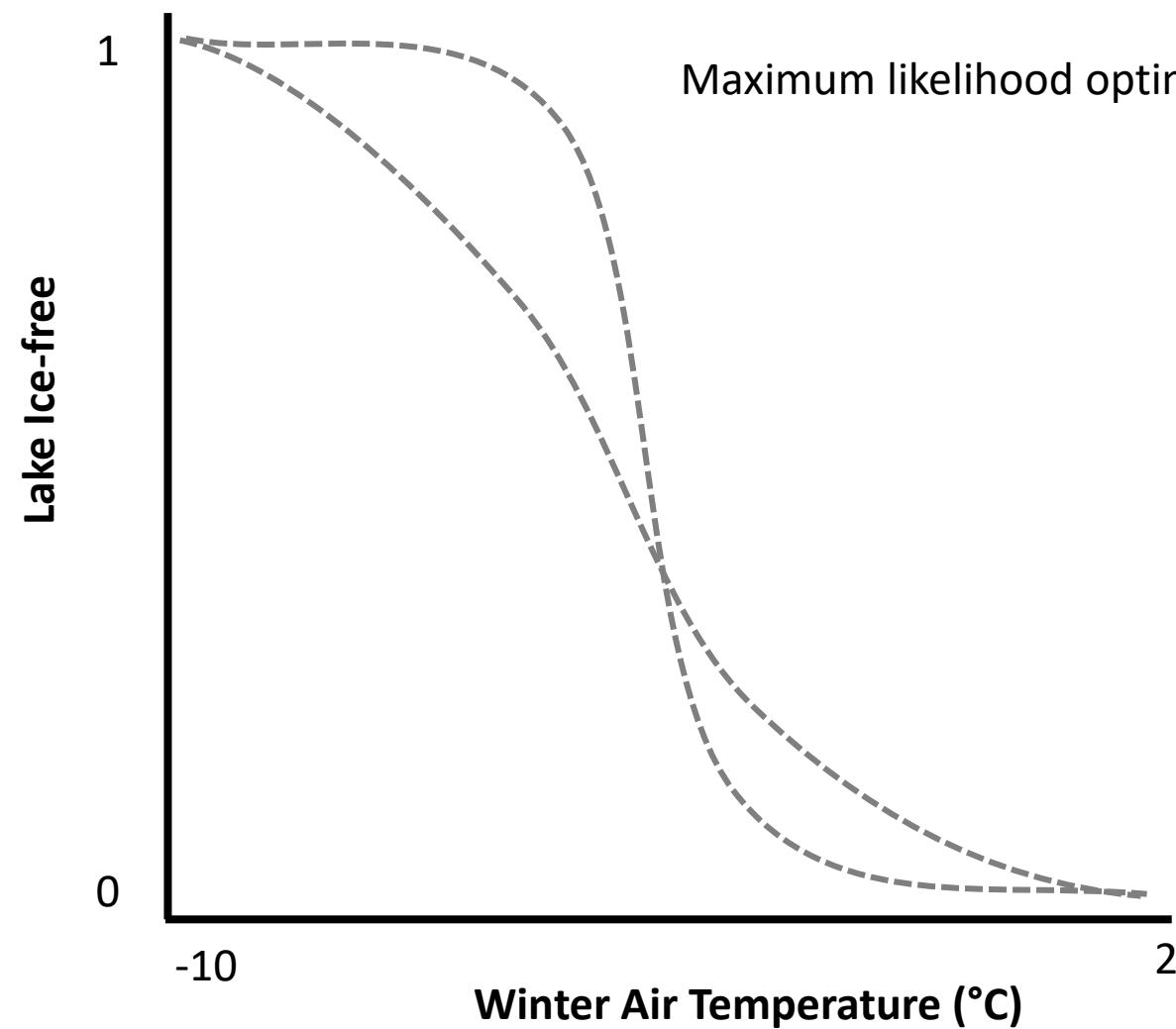
# 1) Fitting the model – maximum likelihood



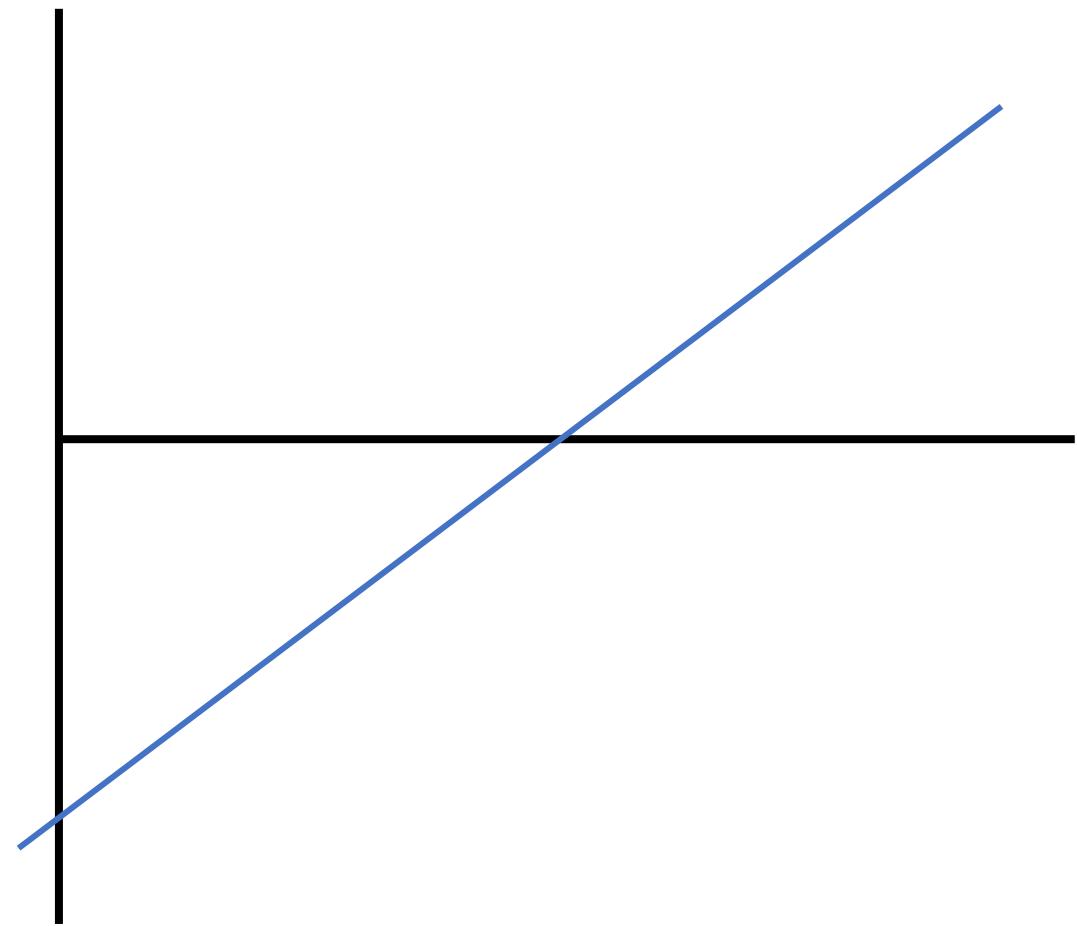
# 1) Fitting the model – maximum likelihood



# 1) Fitting the model – maximum likelihood



Maximum likelihood optimizes quickly

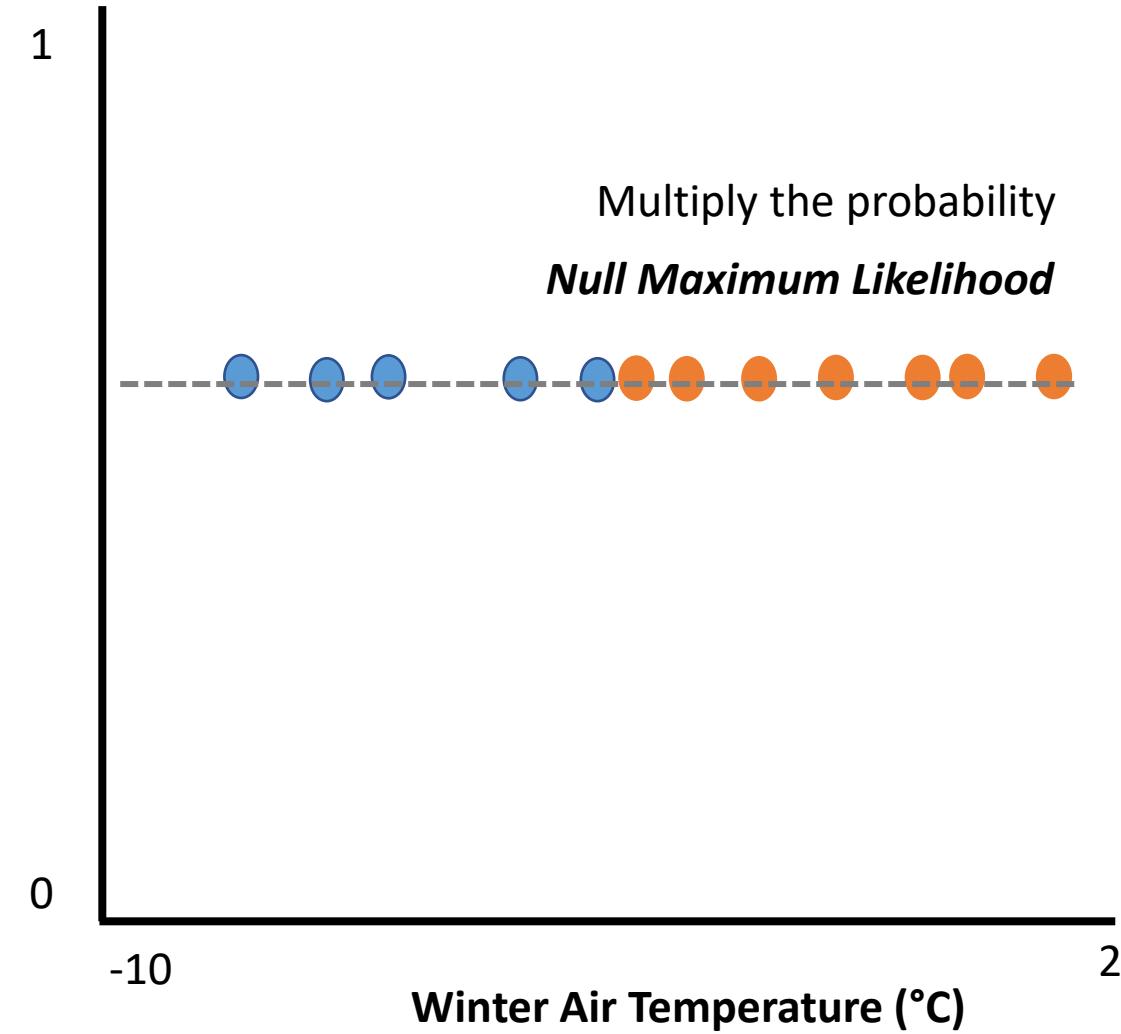
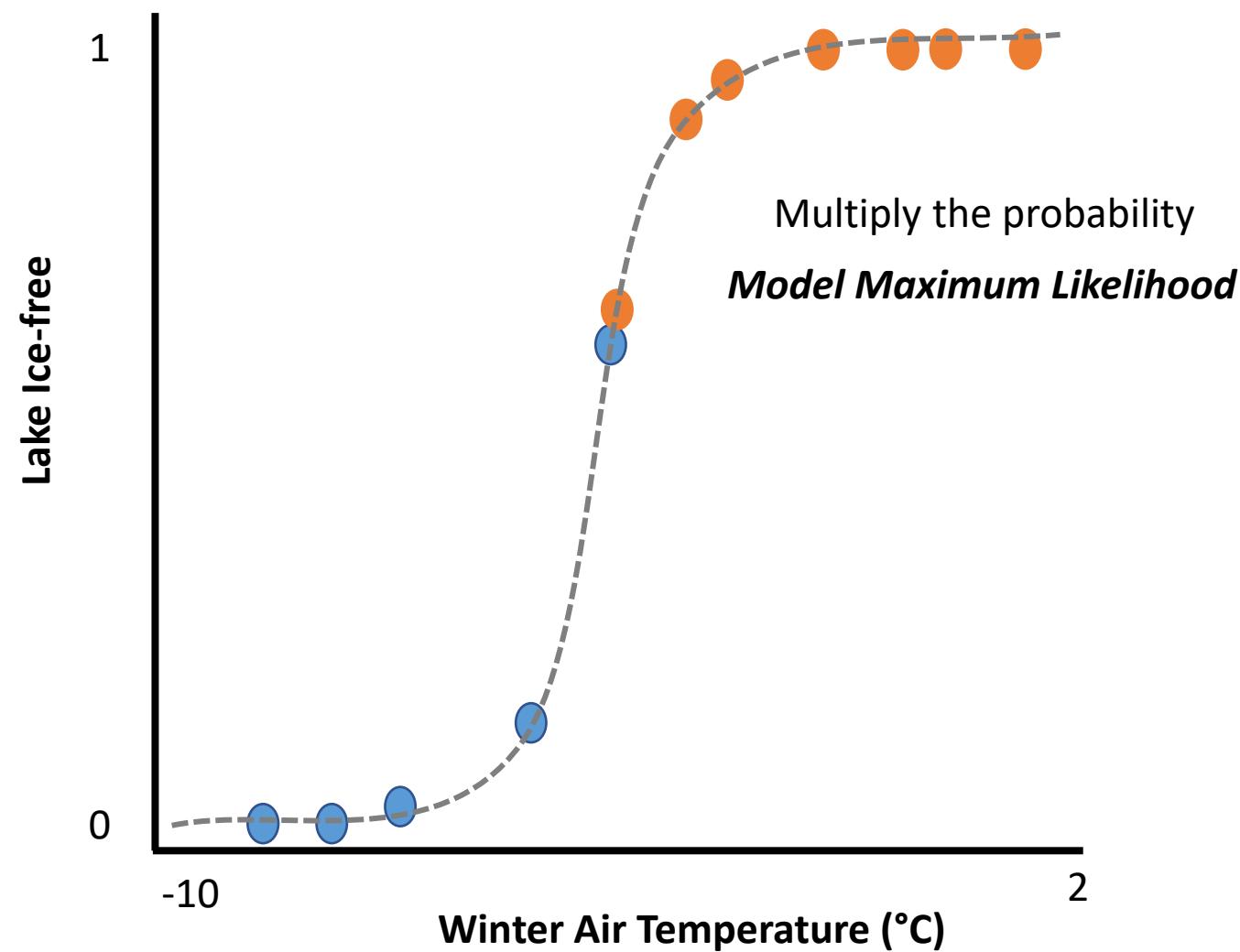


## 2) The outputs – GLMs in R

### 3) Calculating a measure of fit

- Typical linear regression uses R<sup>2</sup>
- No R<sup>2</sup> in logistic regression because no residuals
- Instead use Psuedo R<sup>2</sup>
  - McFadden's among the simpler compares to null

### 3) Calculating a measure of fit



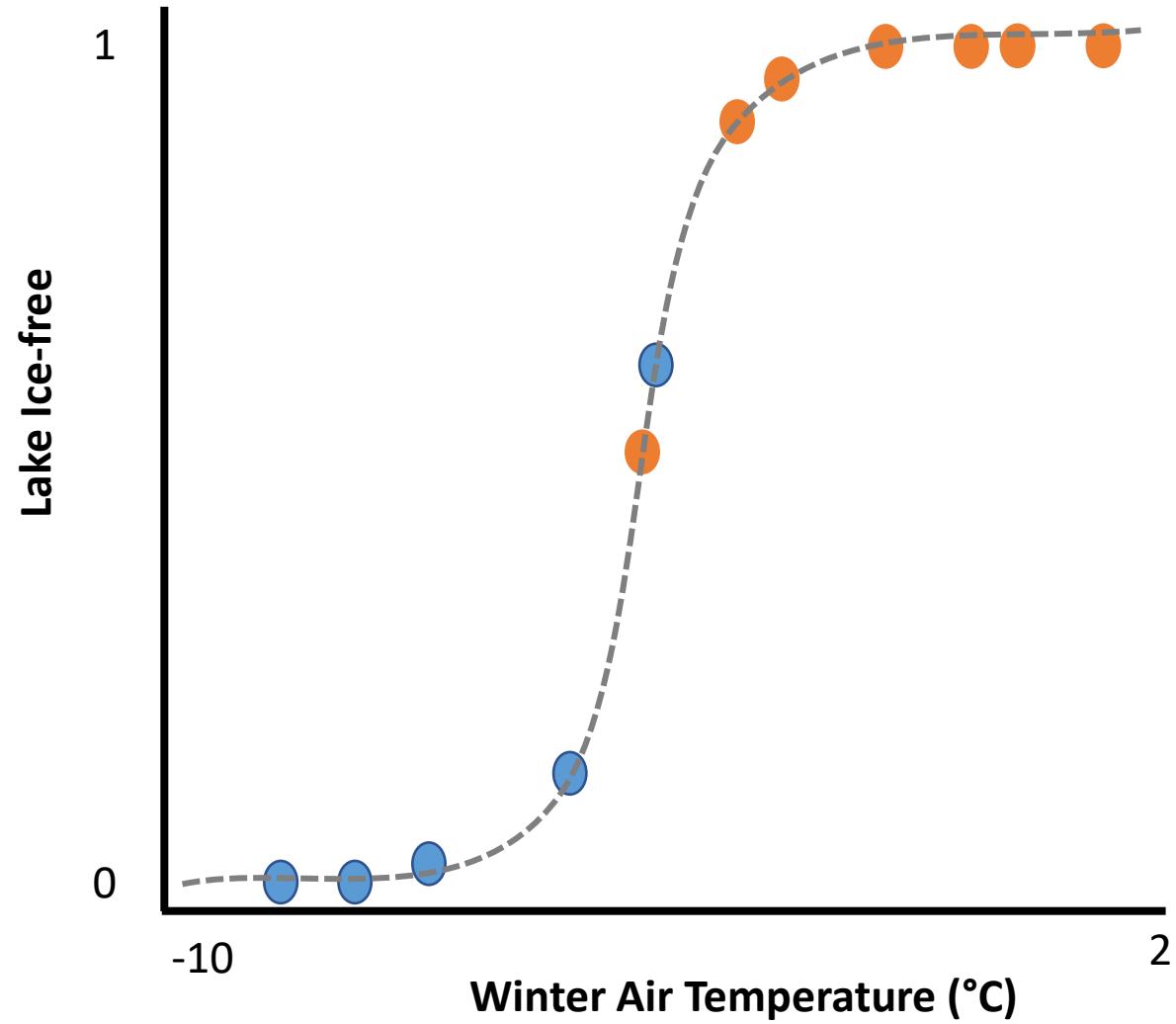
### 3) Calculating a measure of fit

$$\text{PsuedoR}^2 = 1 - \frac{\log(\text{Model})}{\log(\text{Null})}$$

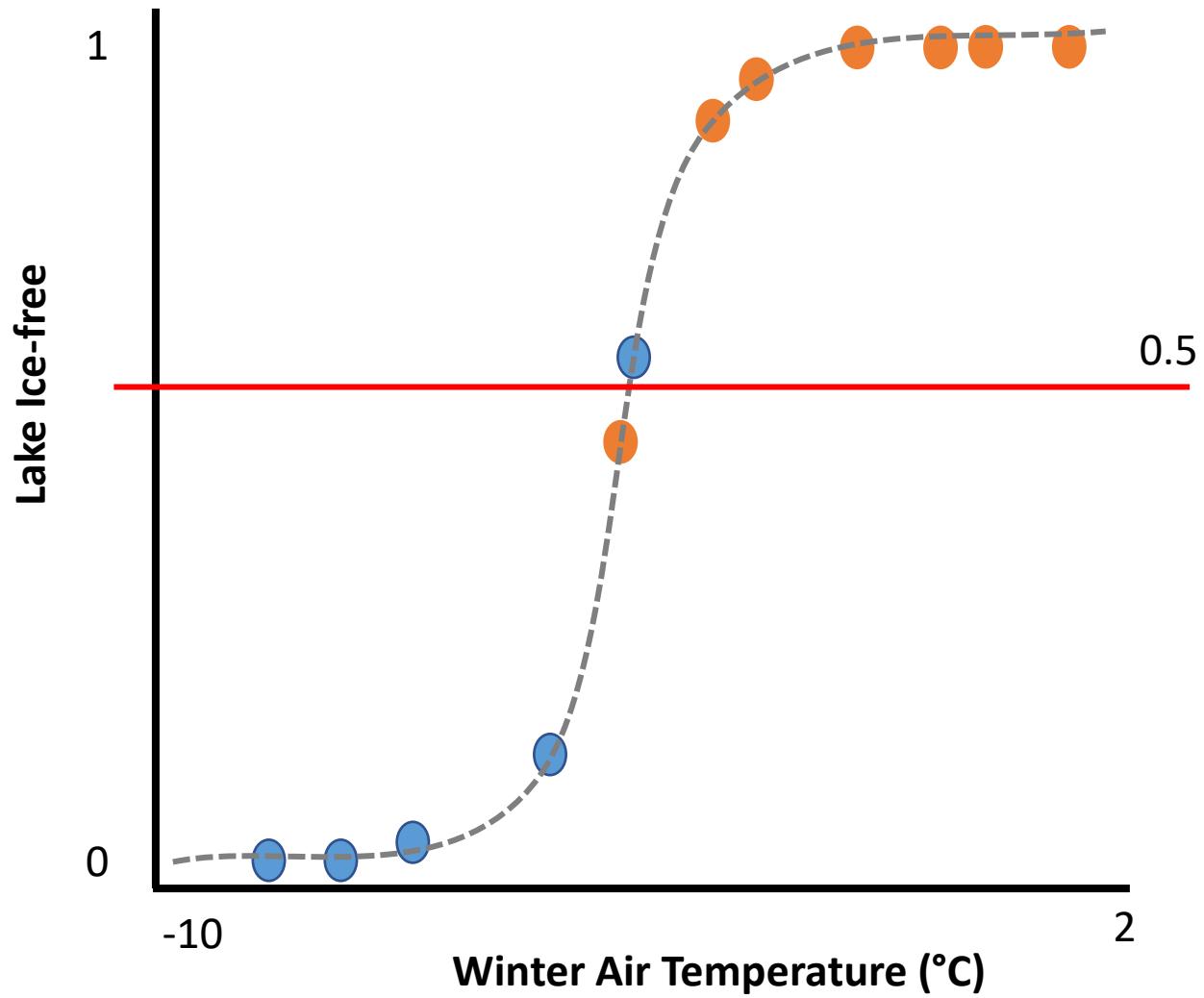
```
m3<- glm(icefree~temperature, family="binomial", data=iceData)
m3null<- glm(icefree~1, family="binomial", data=iceData)
1-logLik(m3)/logLik(m3null)
```

# Prediction

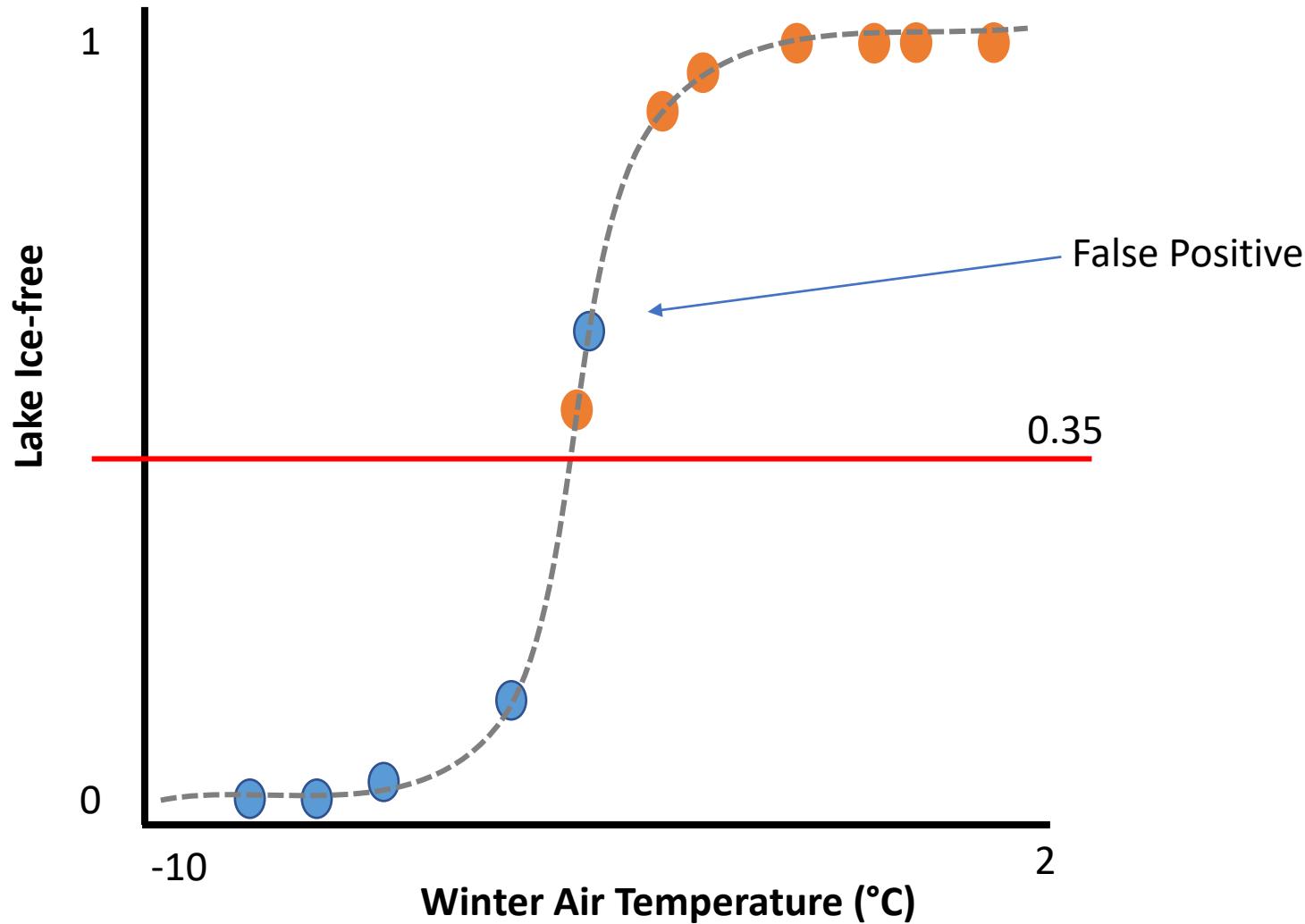
- Models probability
- ..but the outcome is binary
- Need to identify a threshold



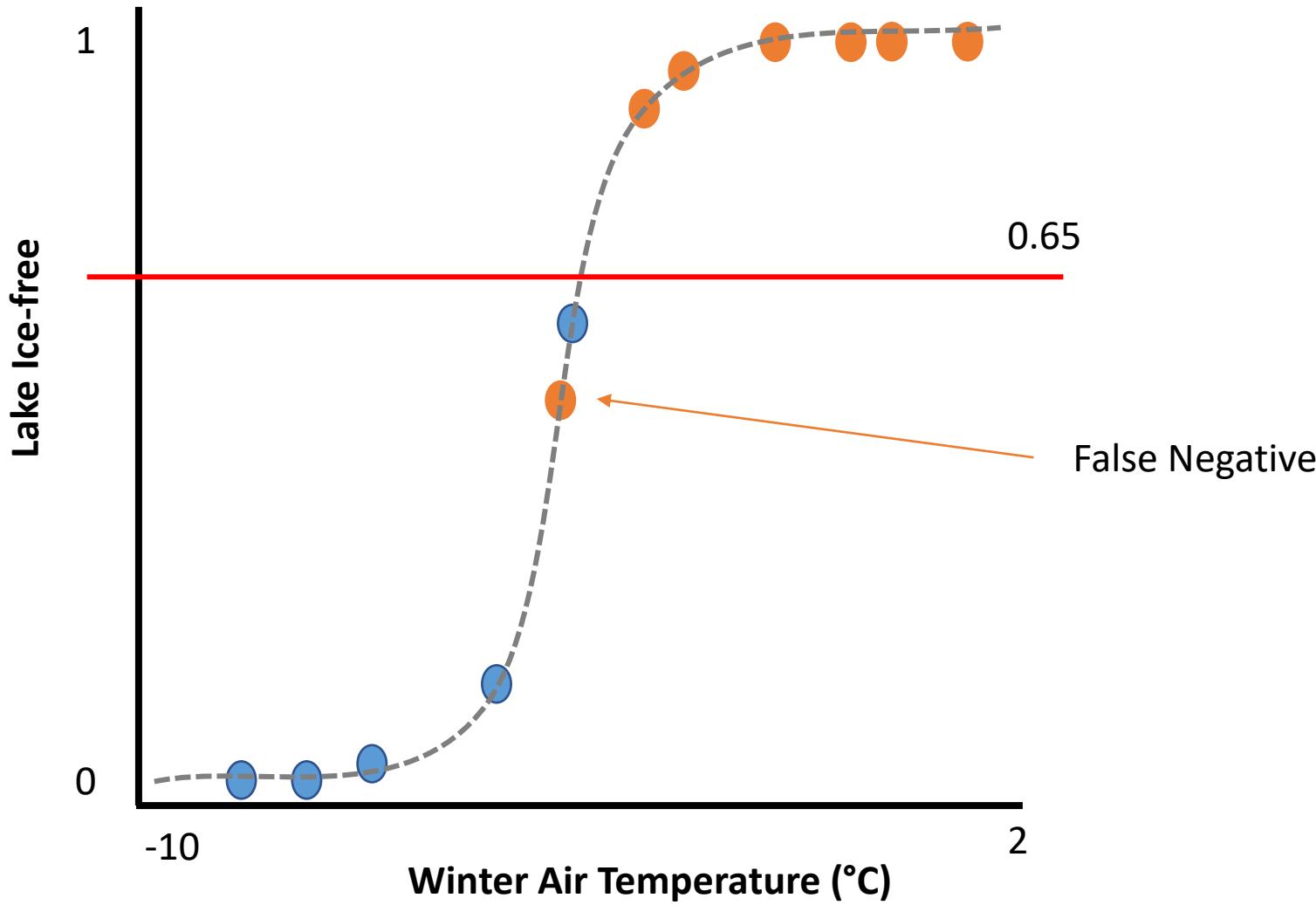
# Prediction



# Prediction



# Prediction

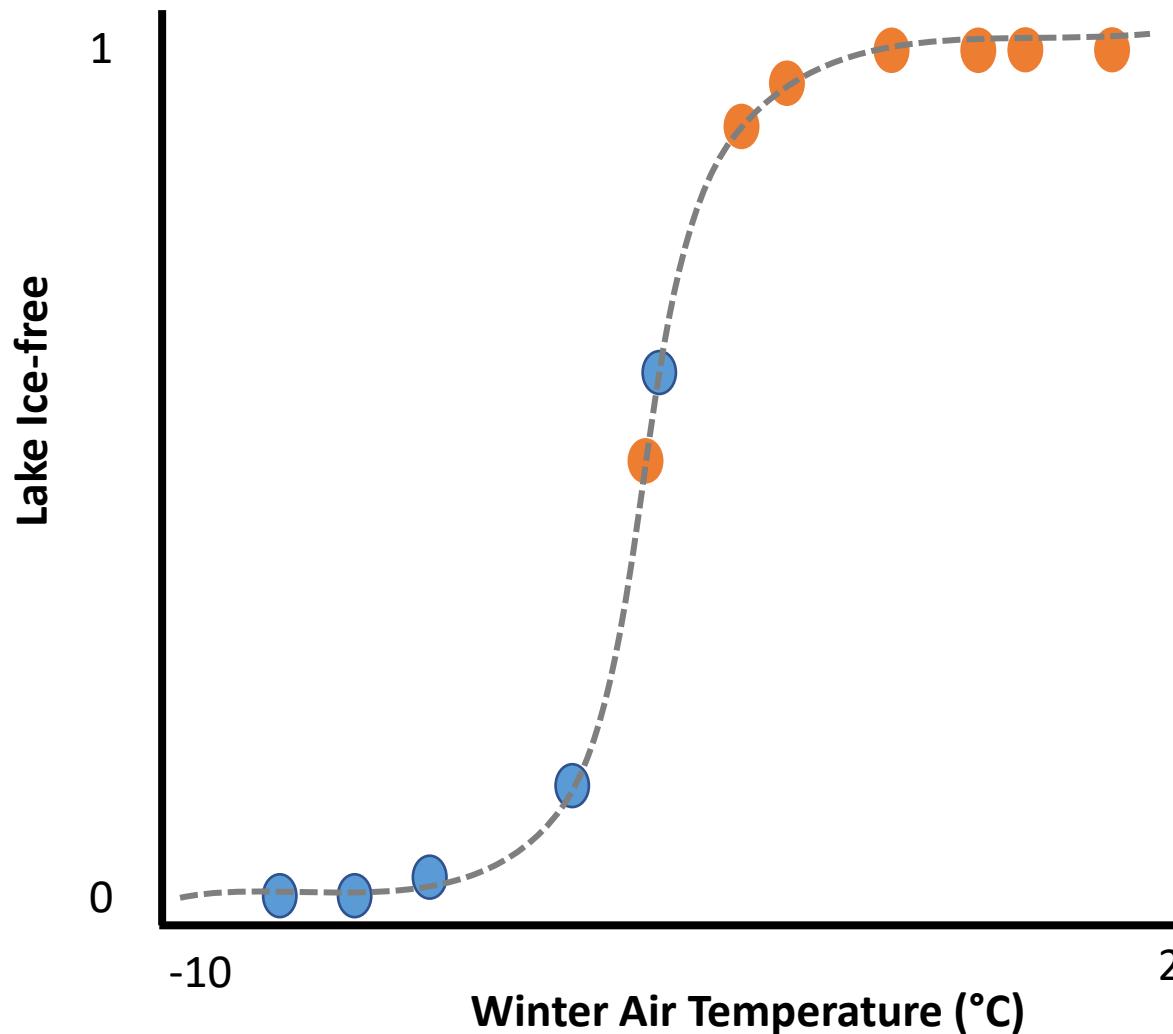


# Question

When is an example where you would like to reduce false negative rate at the expense of false positive?

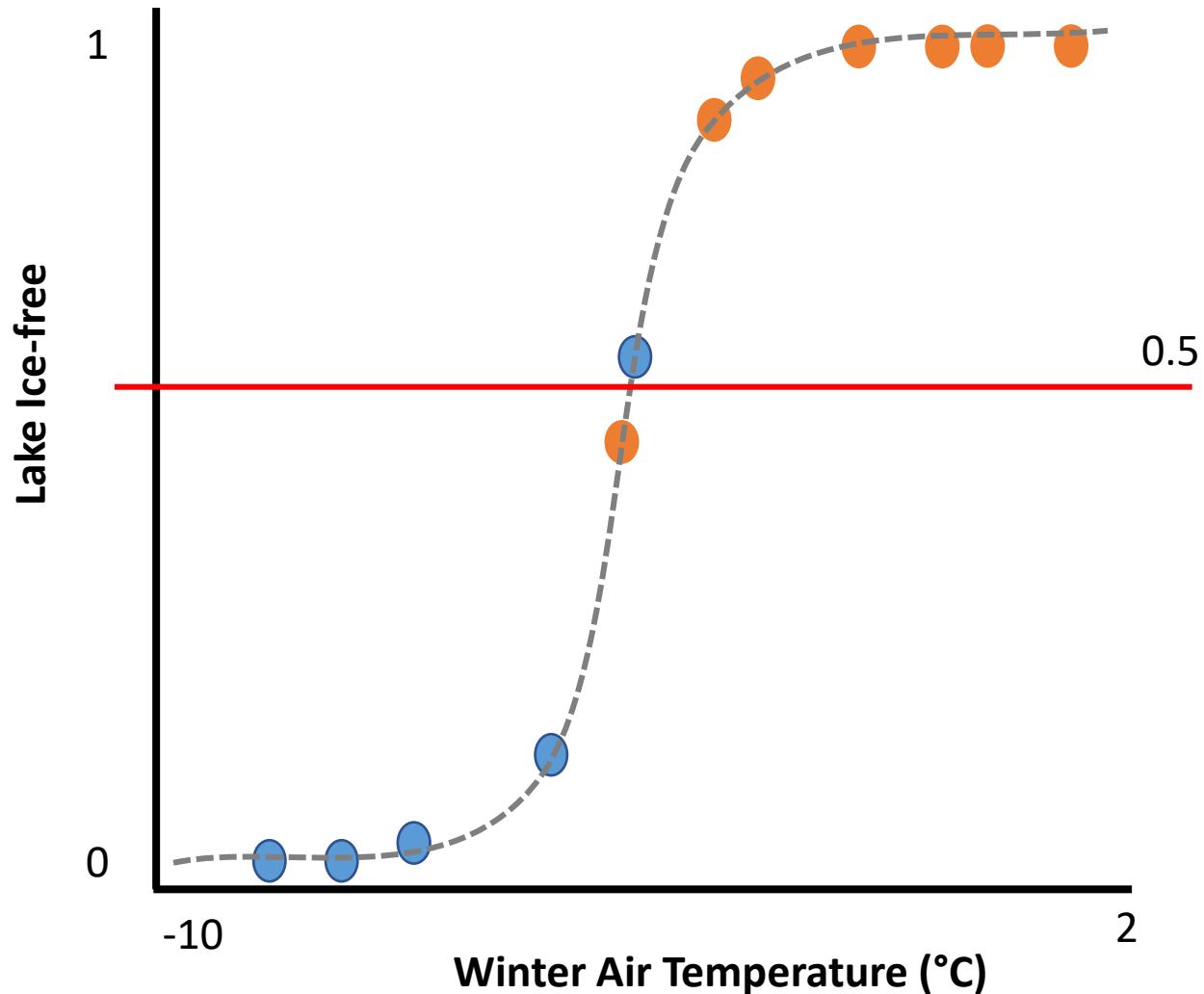
- a) Predicting infection of a disease
- b) Understanding climate effects on lake ice
- c) Predicting the probability of a species occurring
- d) A and C
- e) All the above

# Prediction – confusion matrix



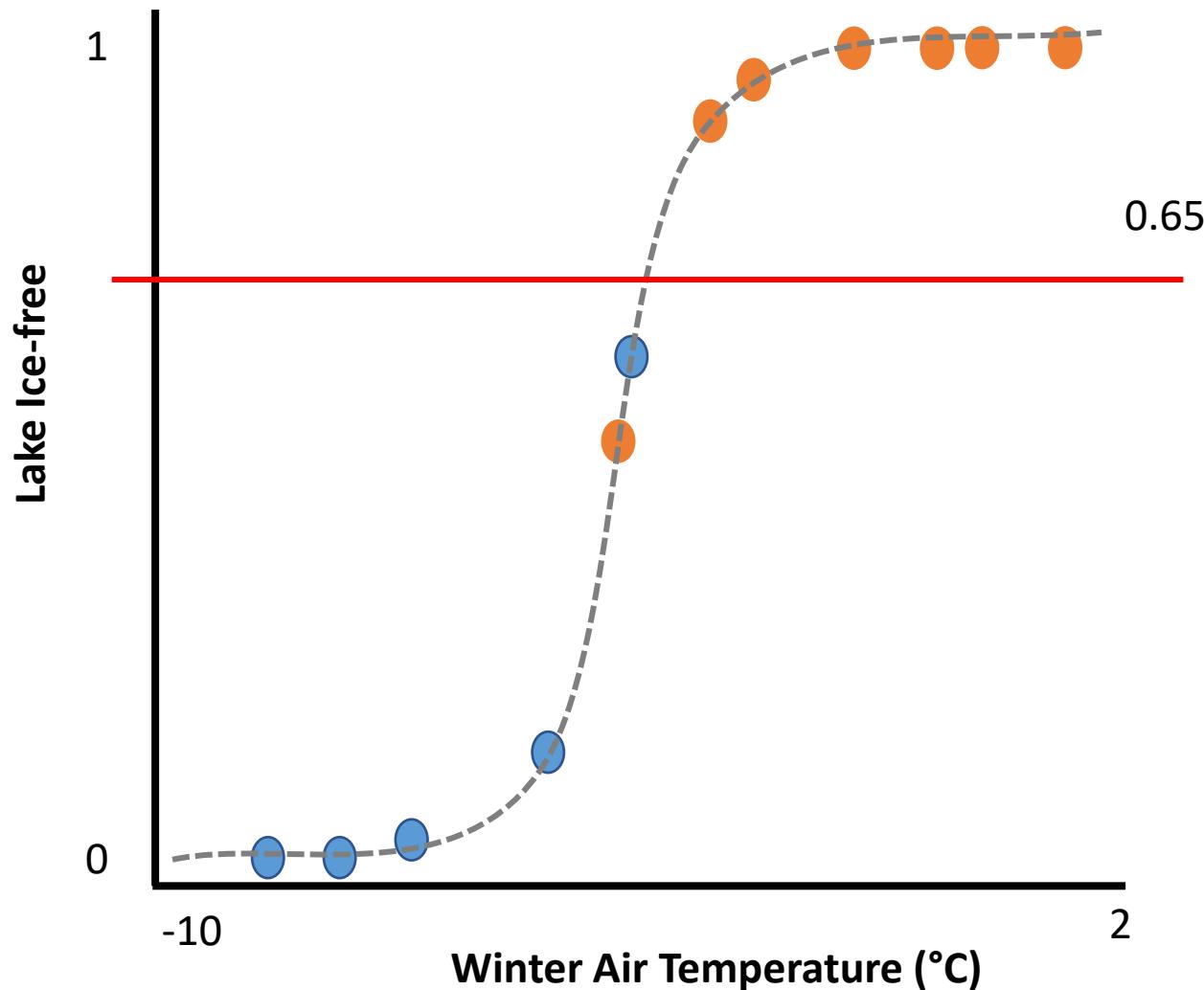
	Actual Positive	Actual Negative
Predicted Positive	True positive	False positive
Predicted Negative	False negative	True negative

# Prediction – confusion matrix



	Actual Positive	Actual Negative
Predicted Positive		
Predicted Negative		

# Prediction – confusion matrix



	Actual Positive	Actual Negative
Predicted Positive		
Predicted Negative		

# Prediction – confusion matrix

HYPOTHESIS TESTING OUTCOMES		R e a l i t y	
		The Null Hypothesis Is True	The Alternative Hypothesis is True
R e s e a r c h	The Null Hypothesis Is True	Accurate $1 - \alpha$ 	Type II Error $\beta$ 
	The Alternative Hypothesis is True	Type I Error $\alpha$ 	Accurate $1 - \beta$ 

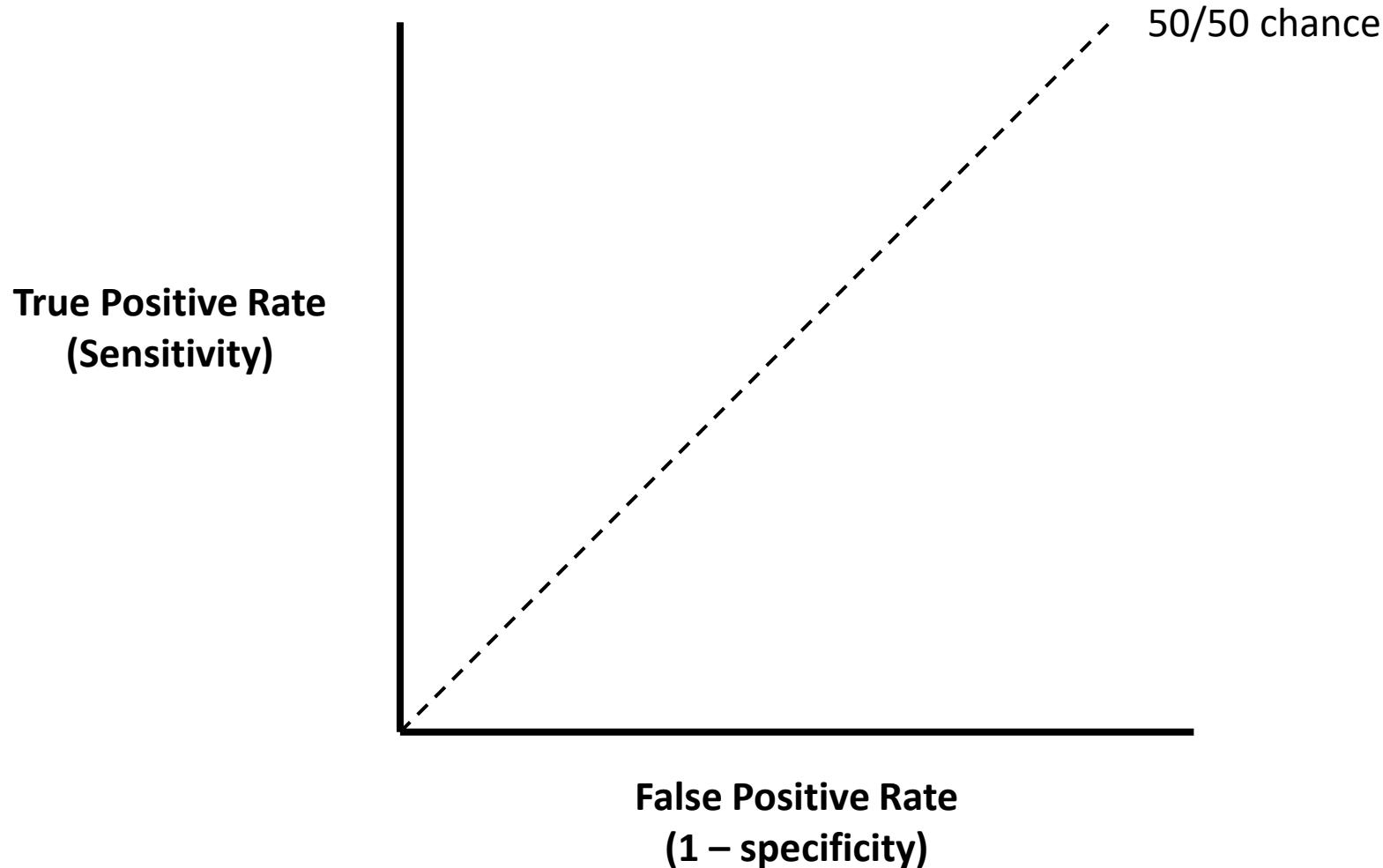
# Optimizing the threshold

Many different methods

Receiver operator characteristic (ROC) often used to pick threshold

Area under the curve (AUC) of ROC used to compare models

# Optimizing the threshold - ROC



# Optimizing the threshold - ROC

	Actual Positive	Actual Negative
Predicted Positive	True positive	False positive
Predicted Negative	False negative	True negative

$$\text{True positive rate} = \frac{\text{True positives}}{\text{Actual positives}}$$

$$\text{True positive rate} = \frac{\text{True positives}}{\text{True positives} + \text{False negatives}}$$

**Sensitivity**

# Optimizing the threshold - ROC

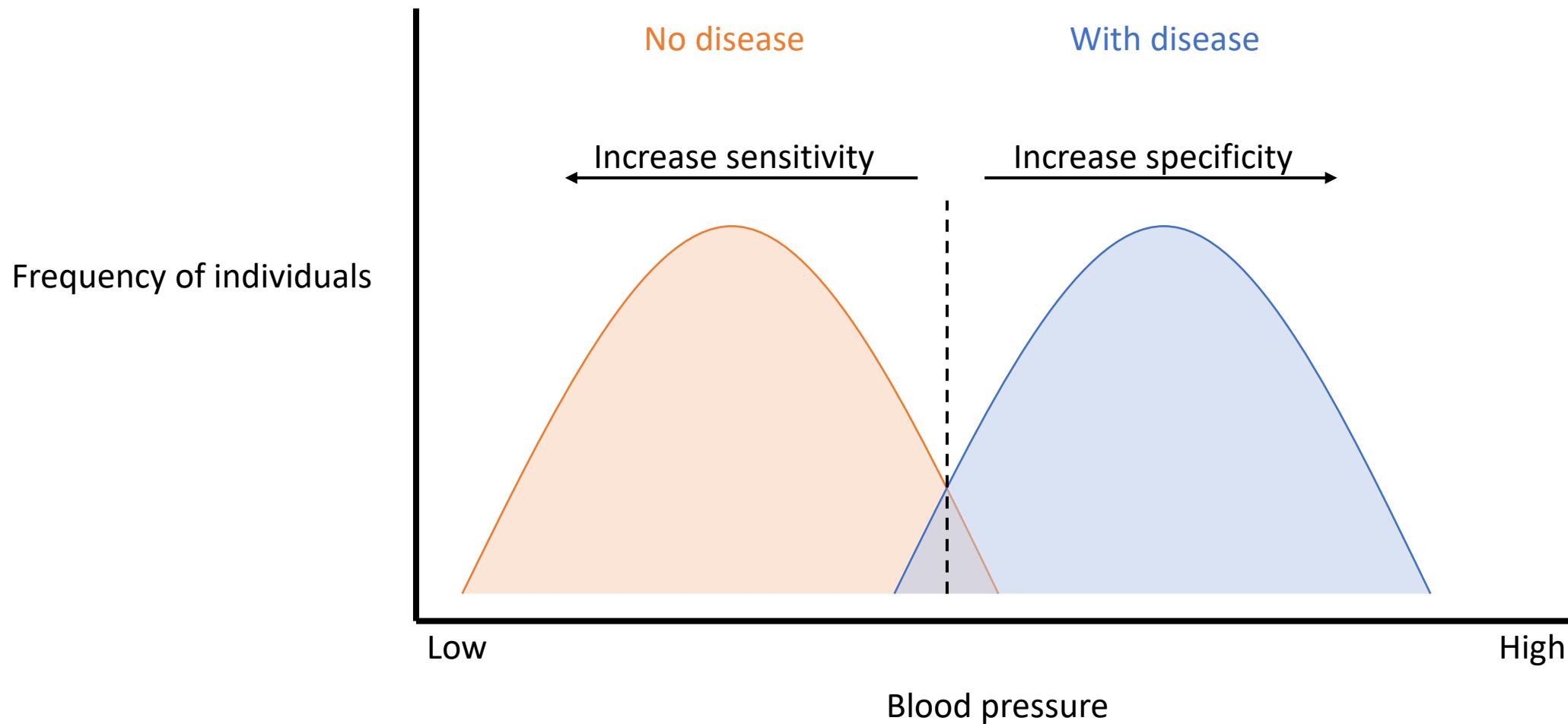
	Actual Positive	Actual Negative
Predicted Positive	True positive	False positive
Predicted Negative	False negative	True negative

$$\text{False positive rate} = \frac{\text{False positives}}{\text{Actual negatives}}$$

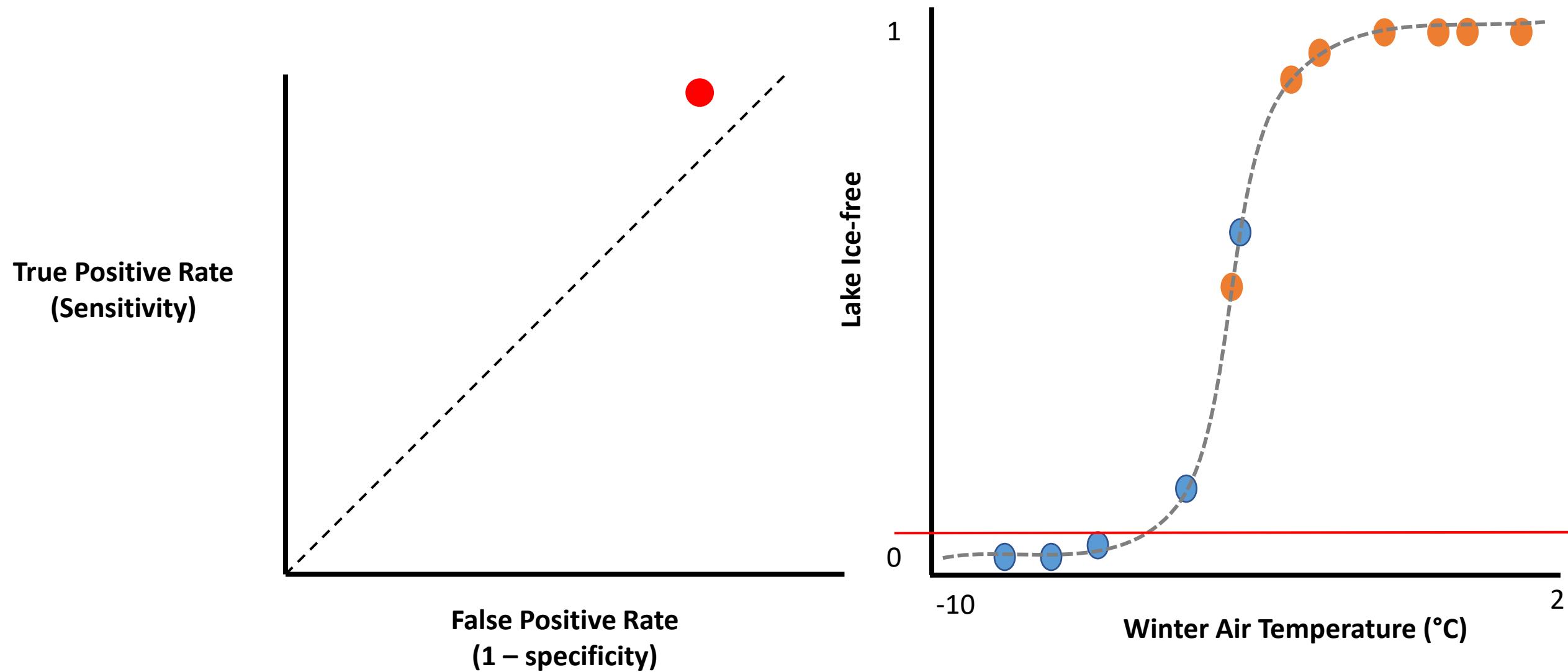
$$\text{False positive rate} = \frac{\text{True positives}}{\text{False positives} + \text{True negatives}}$$

**Specificity**

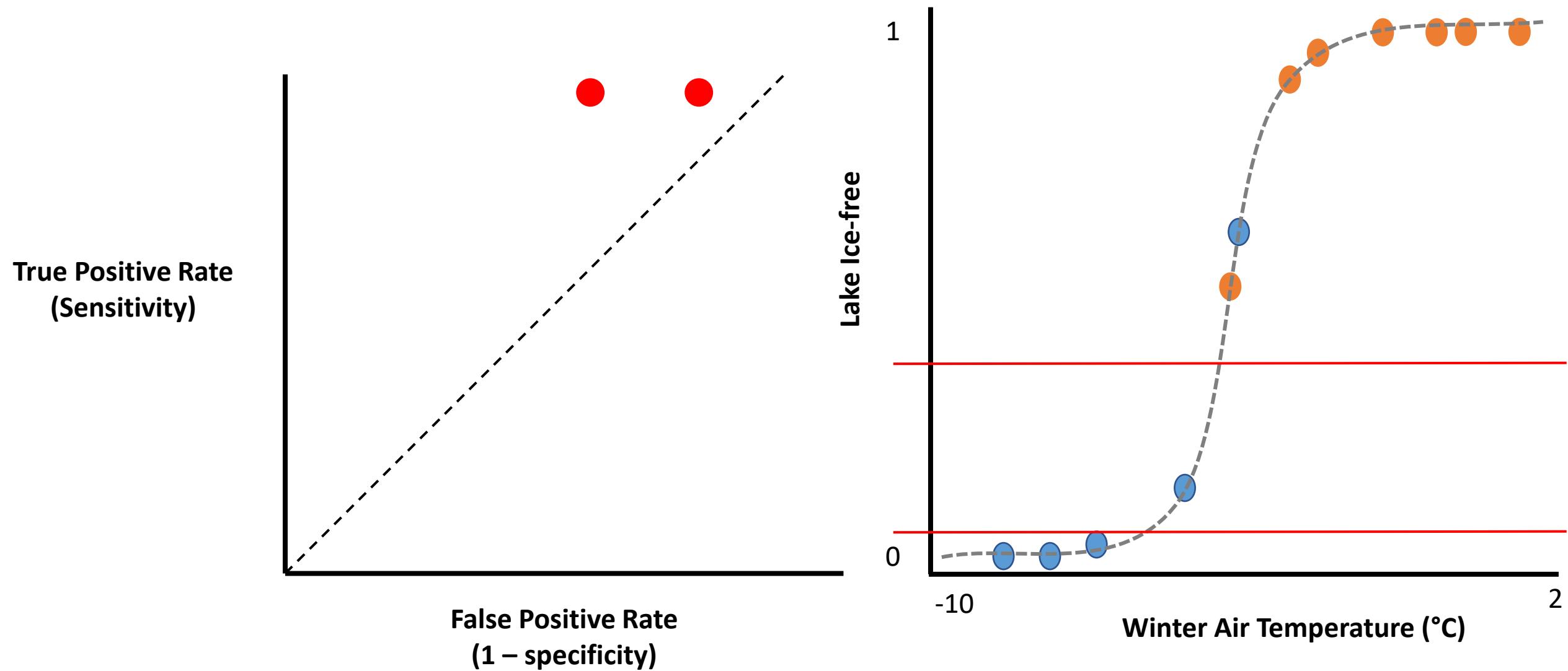
# Sensitivity vs. specificity



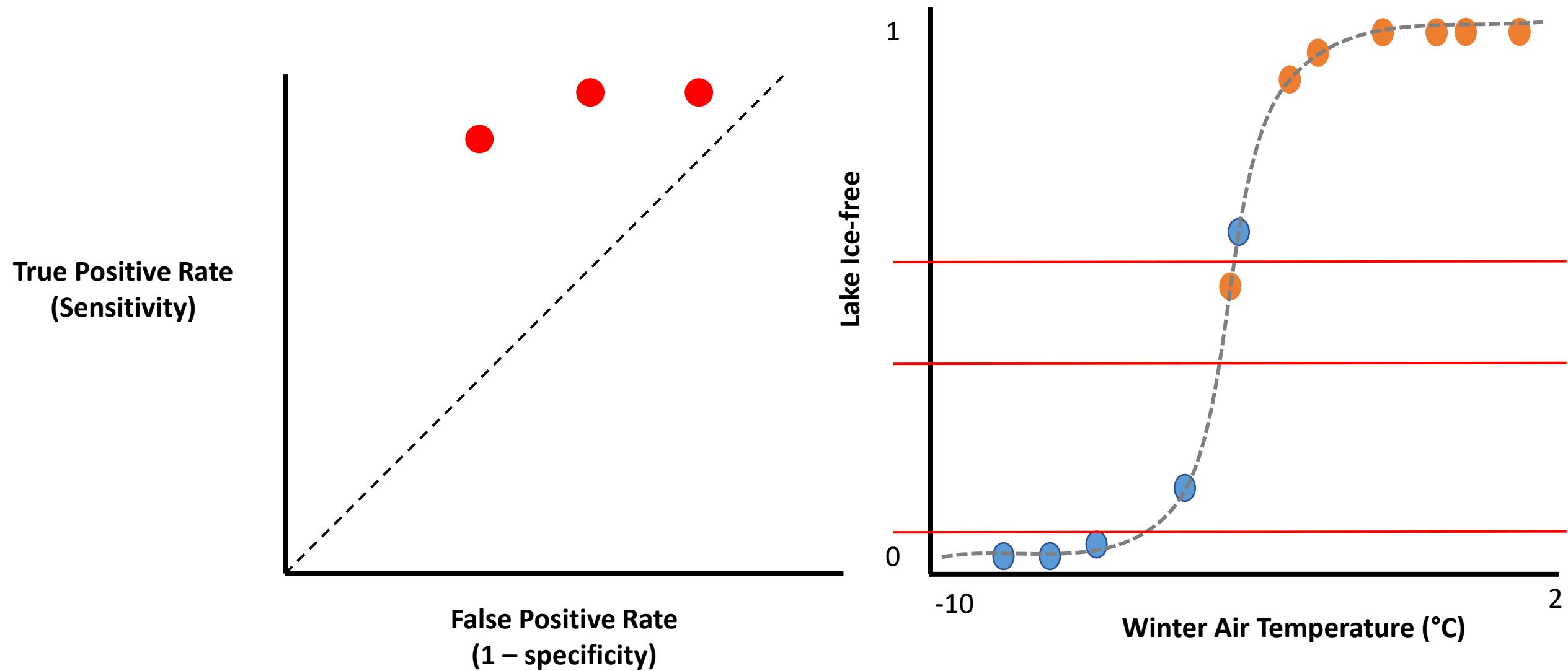
# Optimizing the threshold - ROC



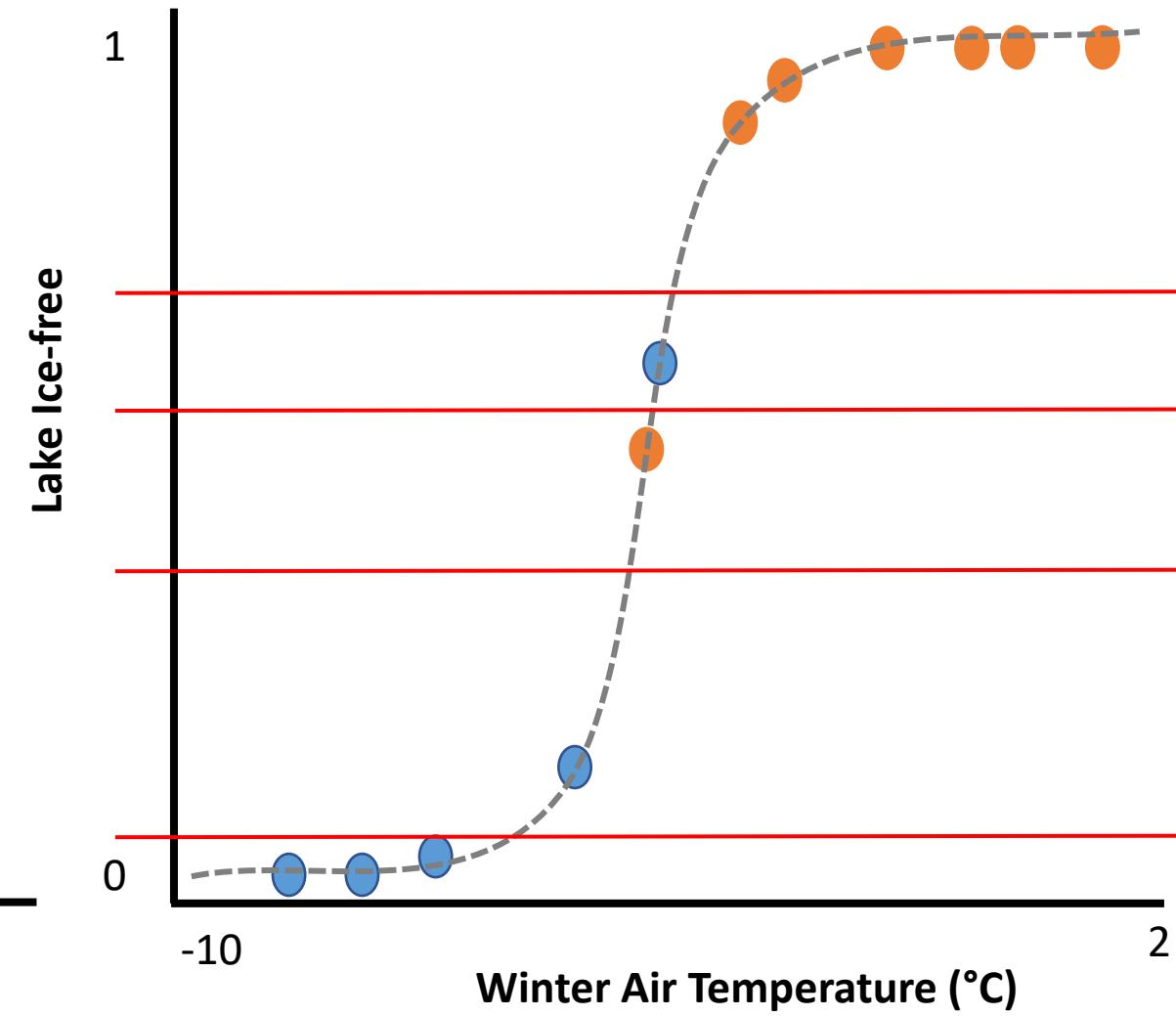
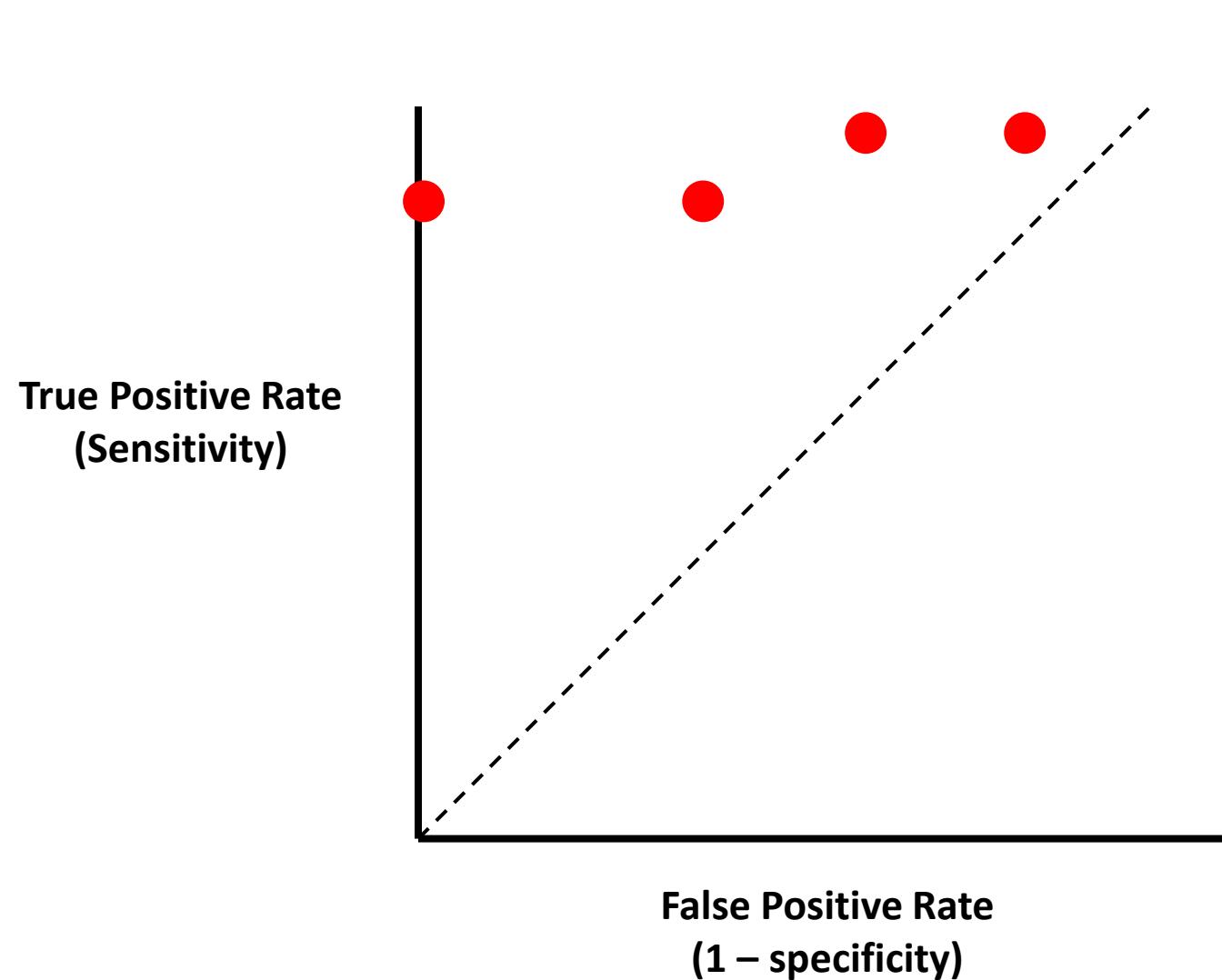
# Optimizing the threshold - ROC



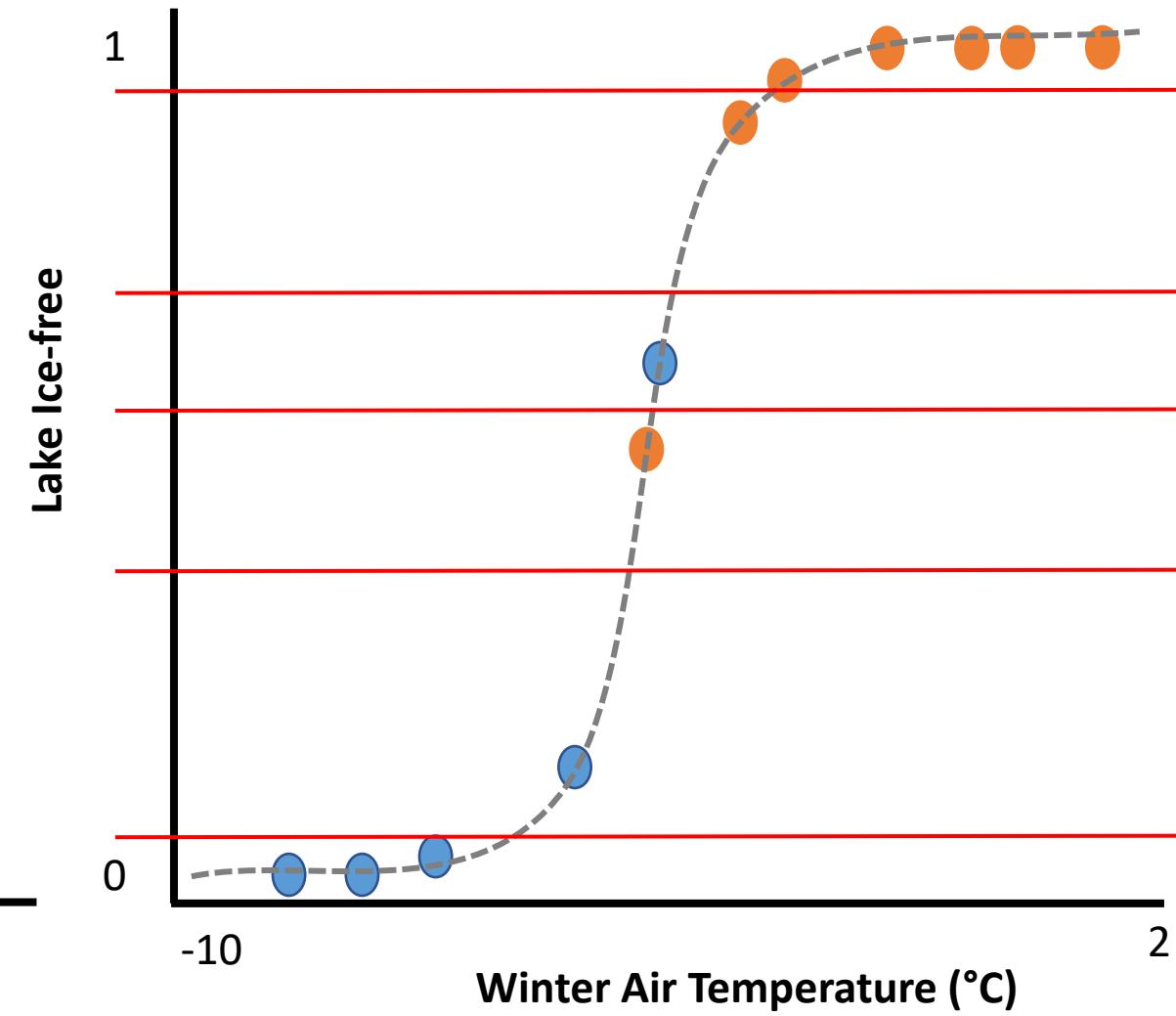
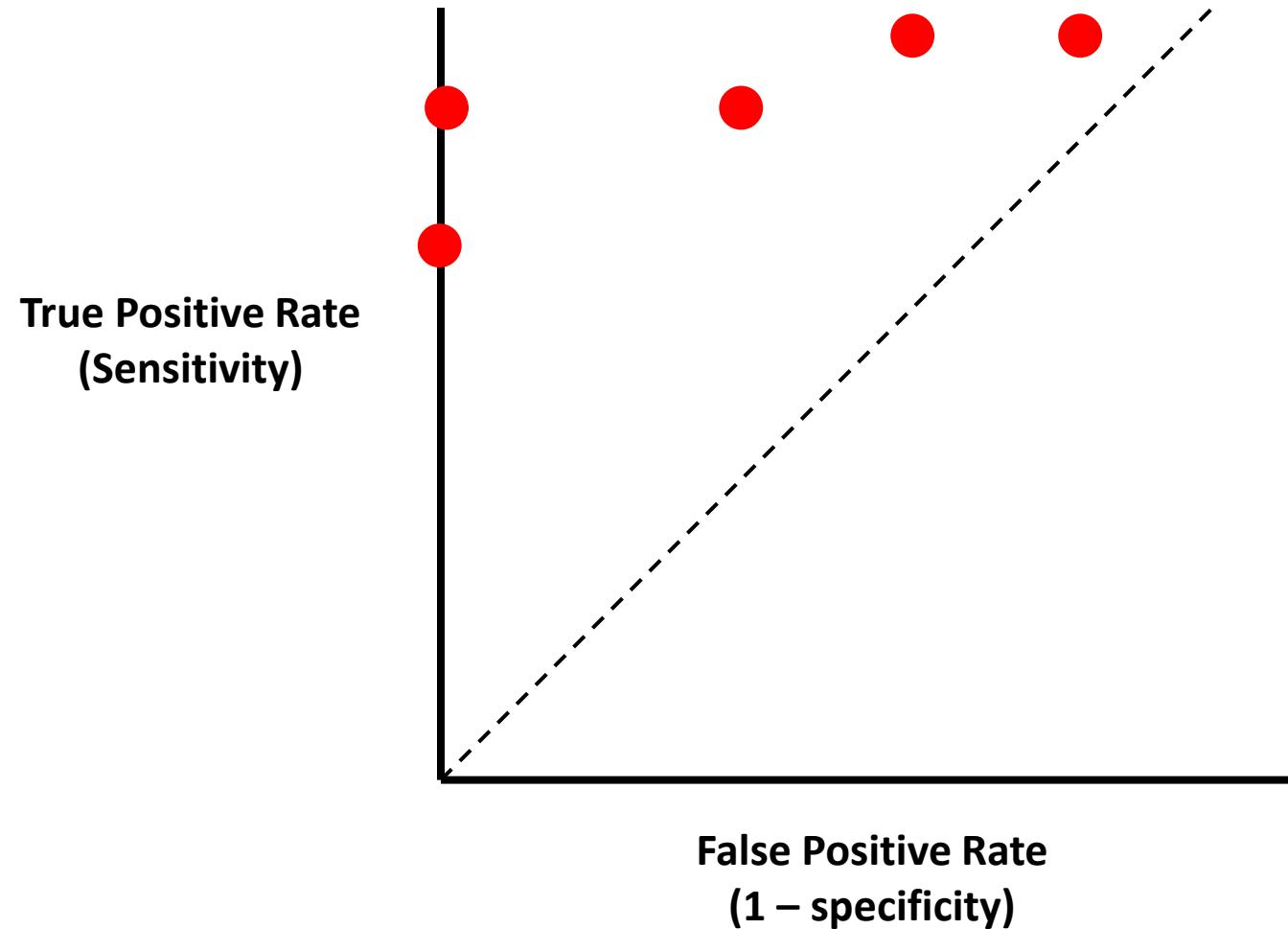
# Optimizing the threshold - ROC



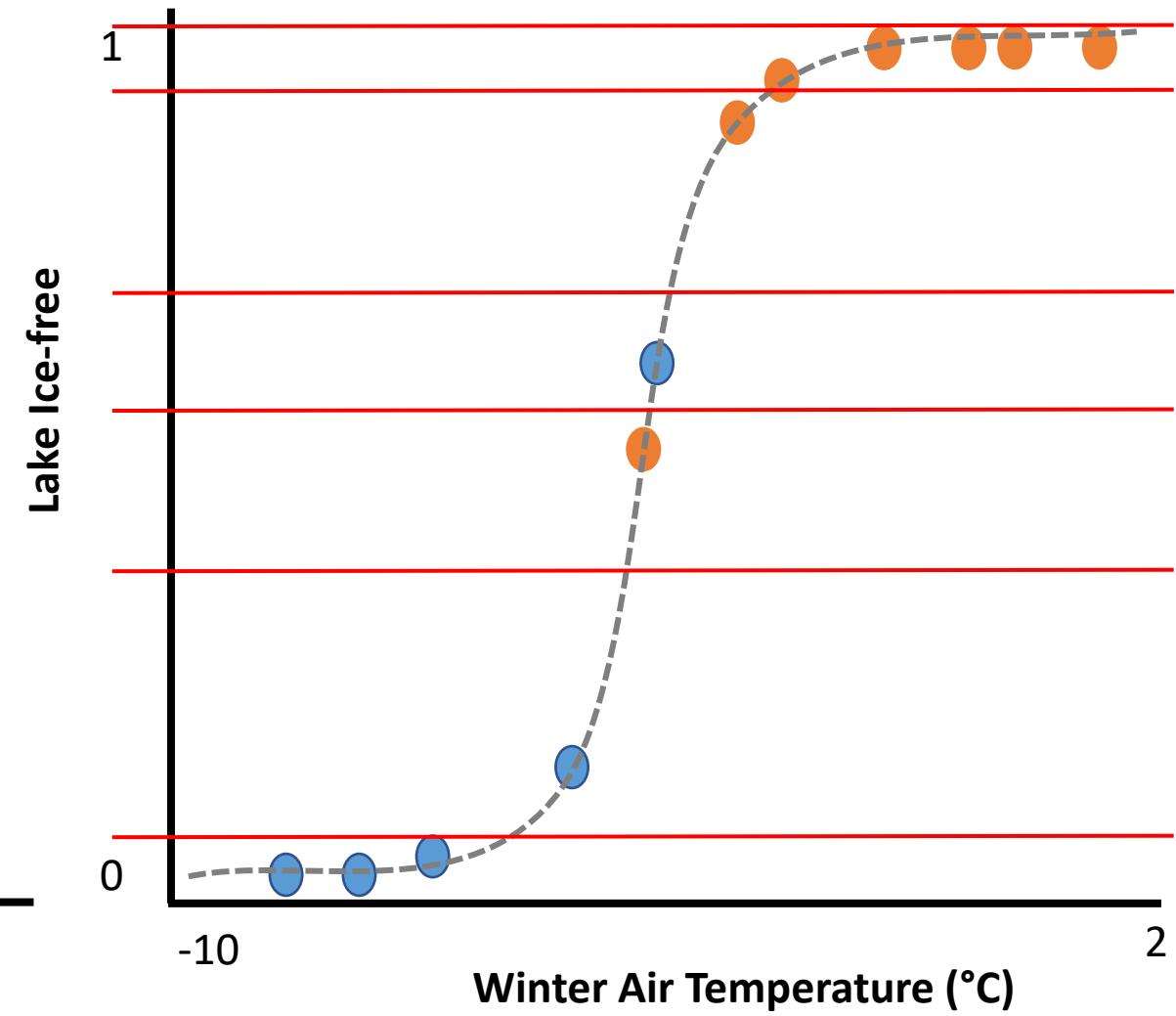
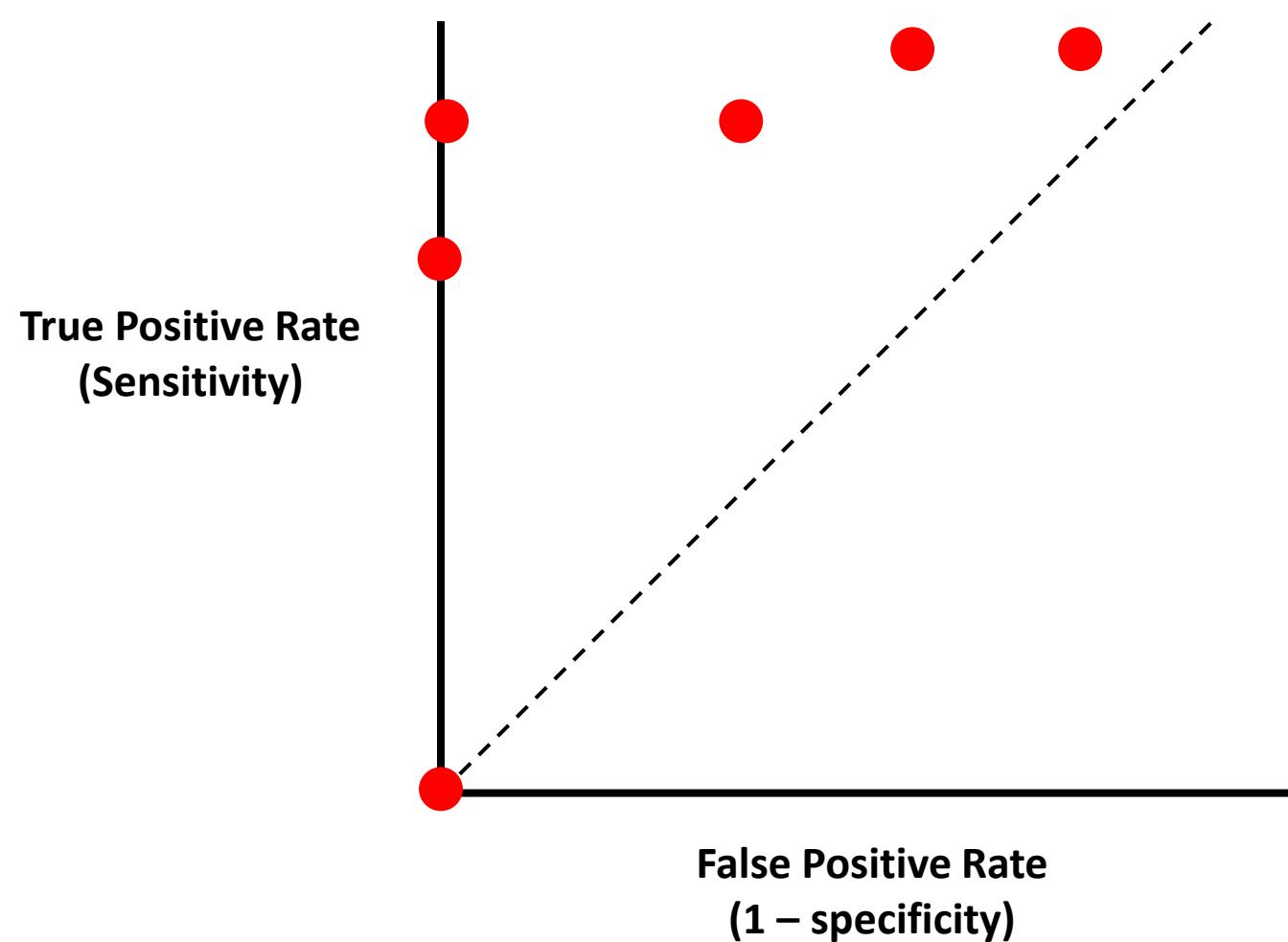
# Optimizing the threshold - ROC



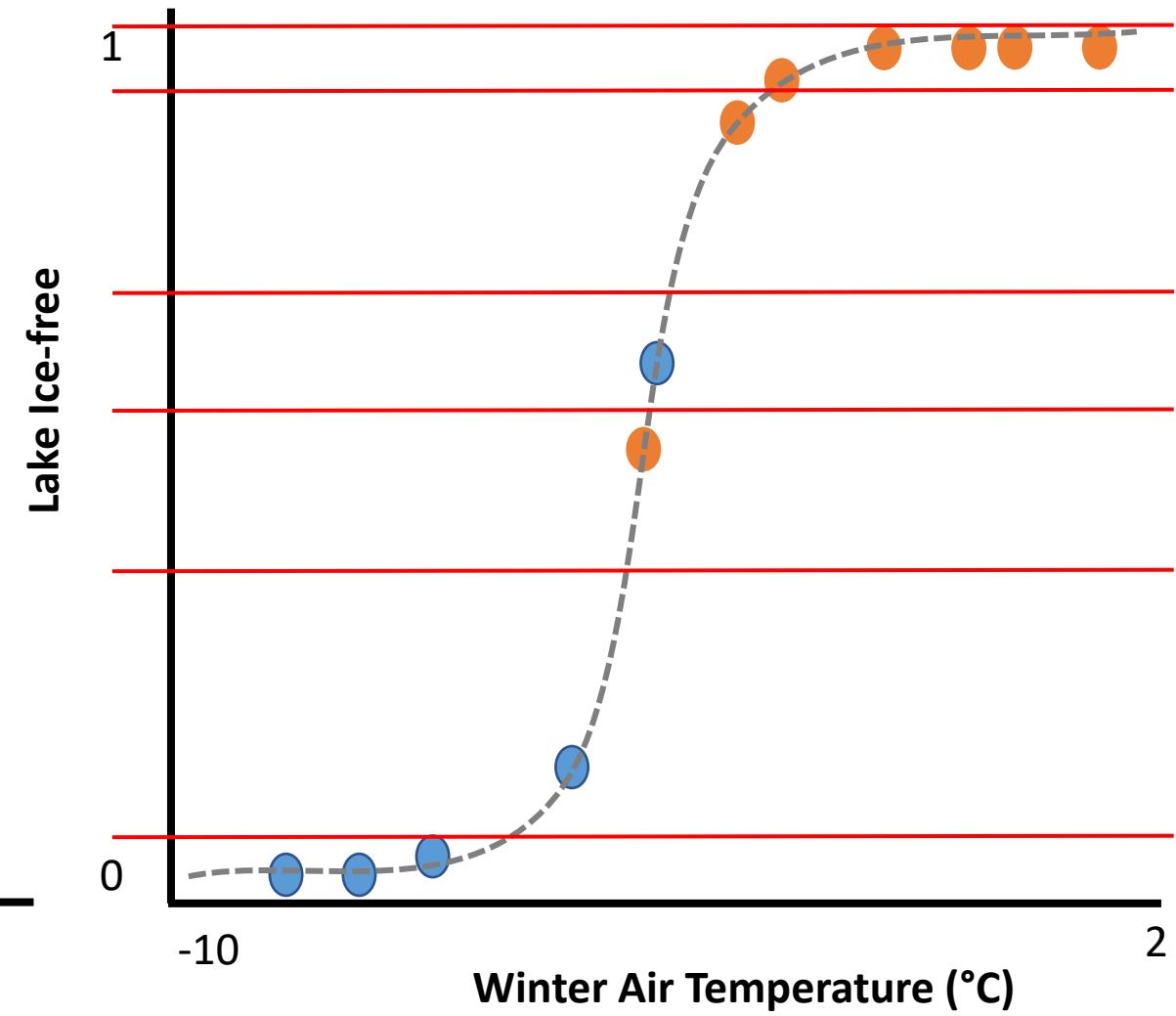
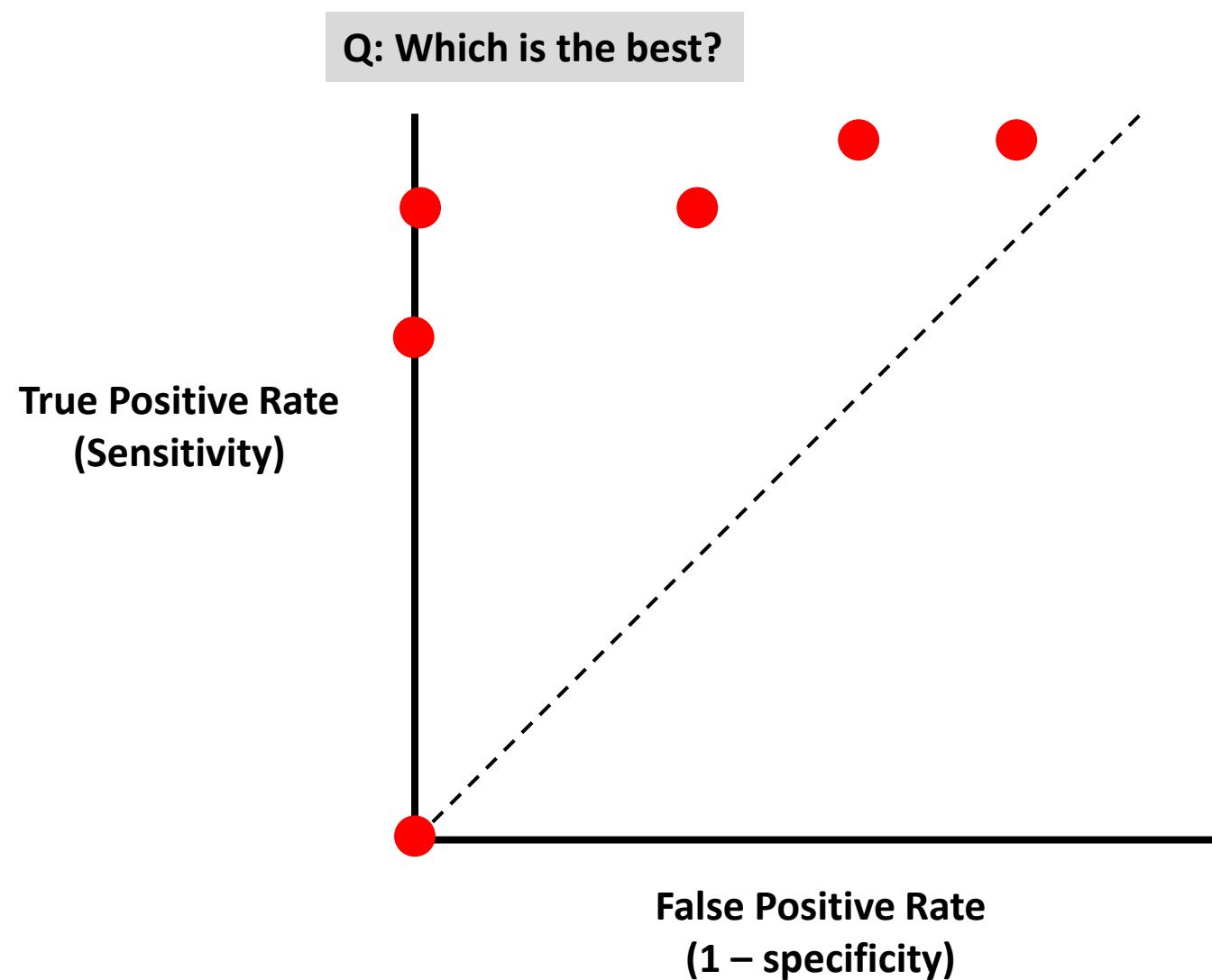
# Optimizing the threshold - ROC



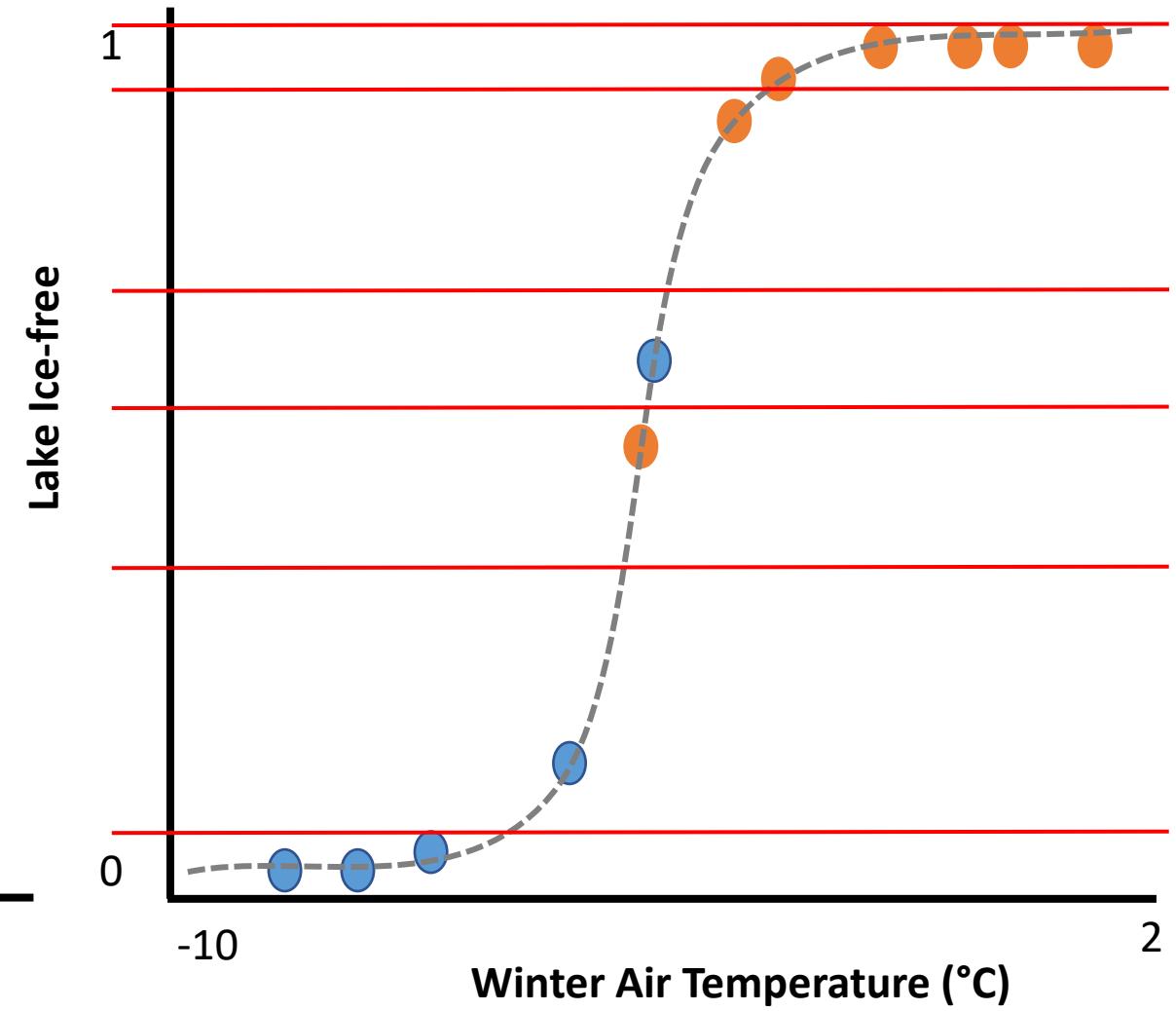
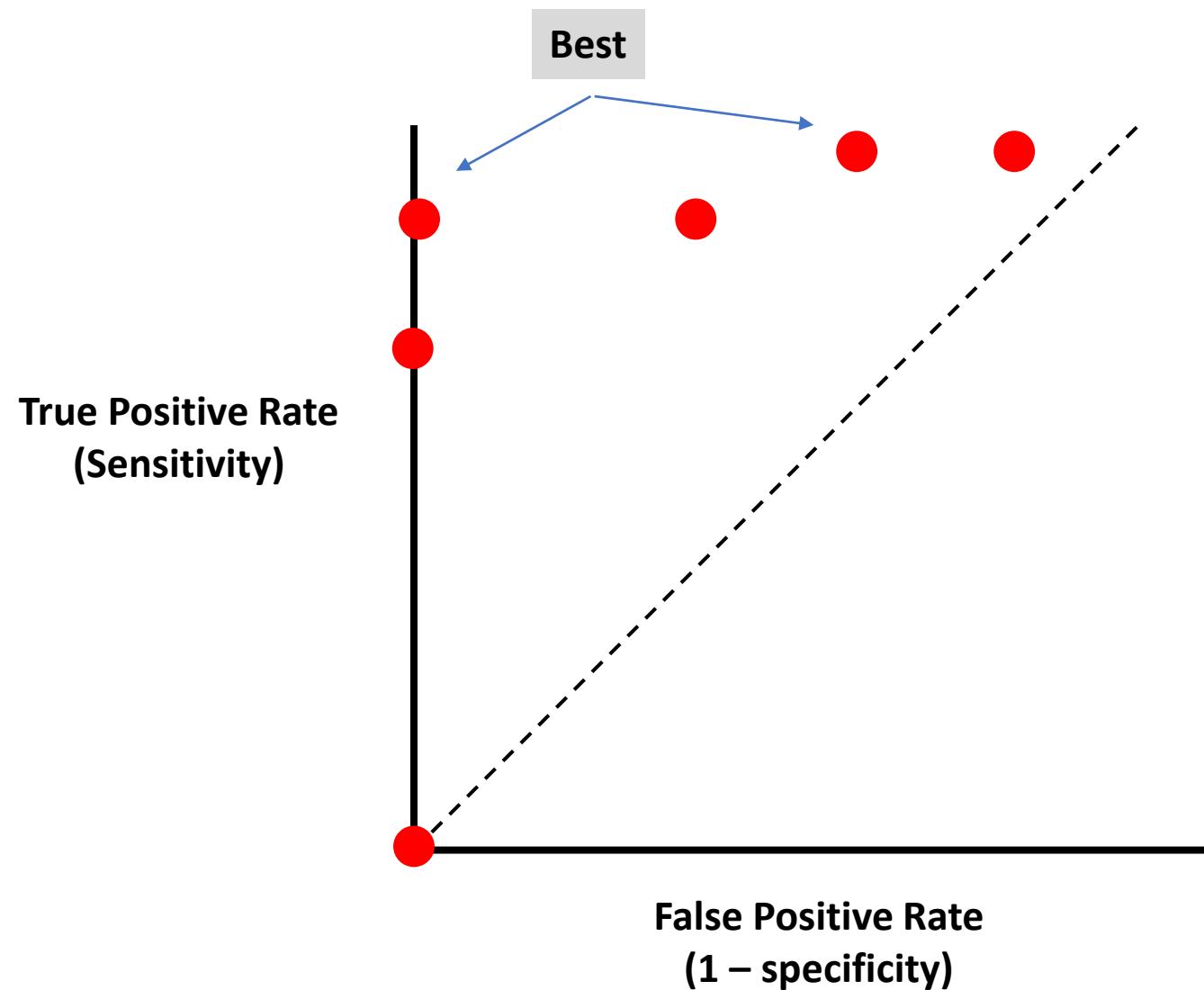
# Optimizing the threshold - ROC



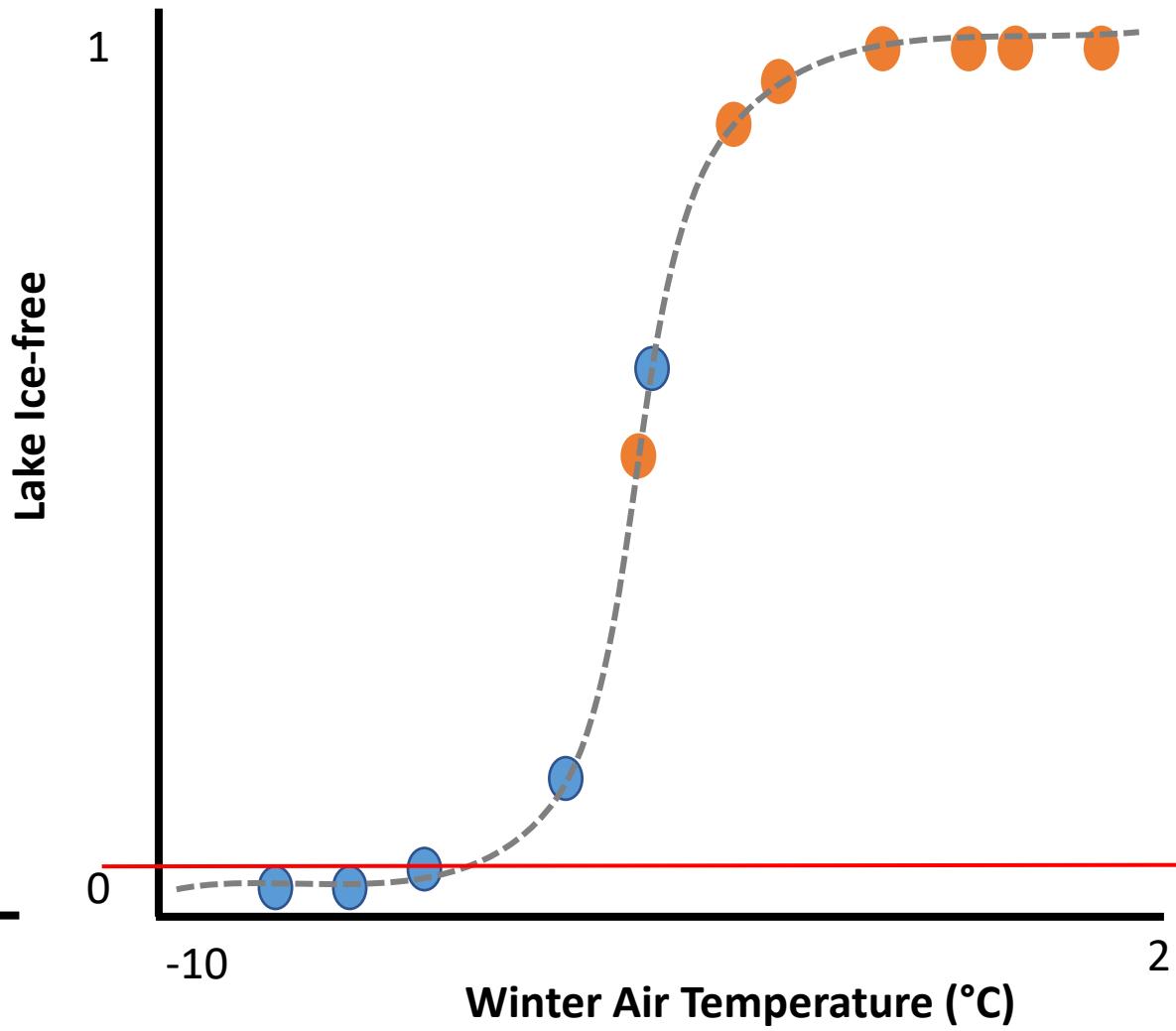
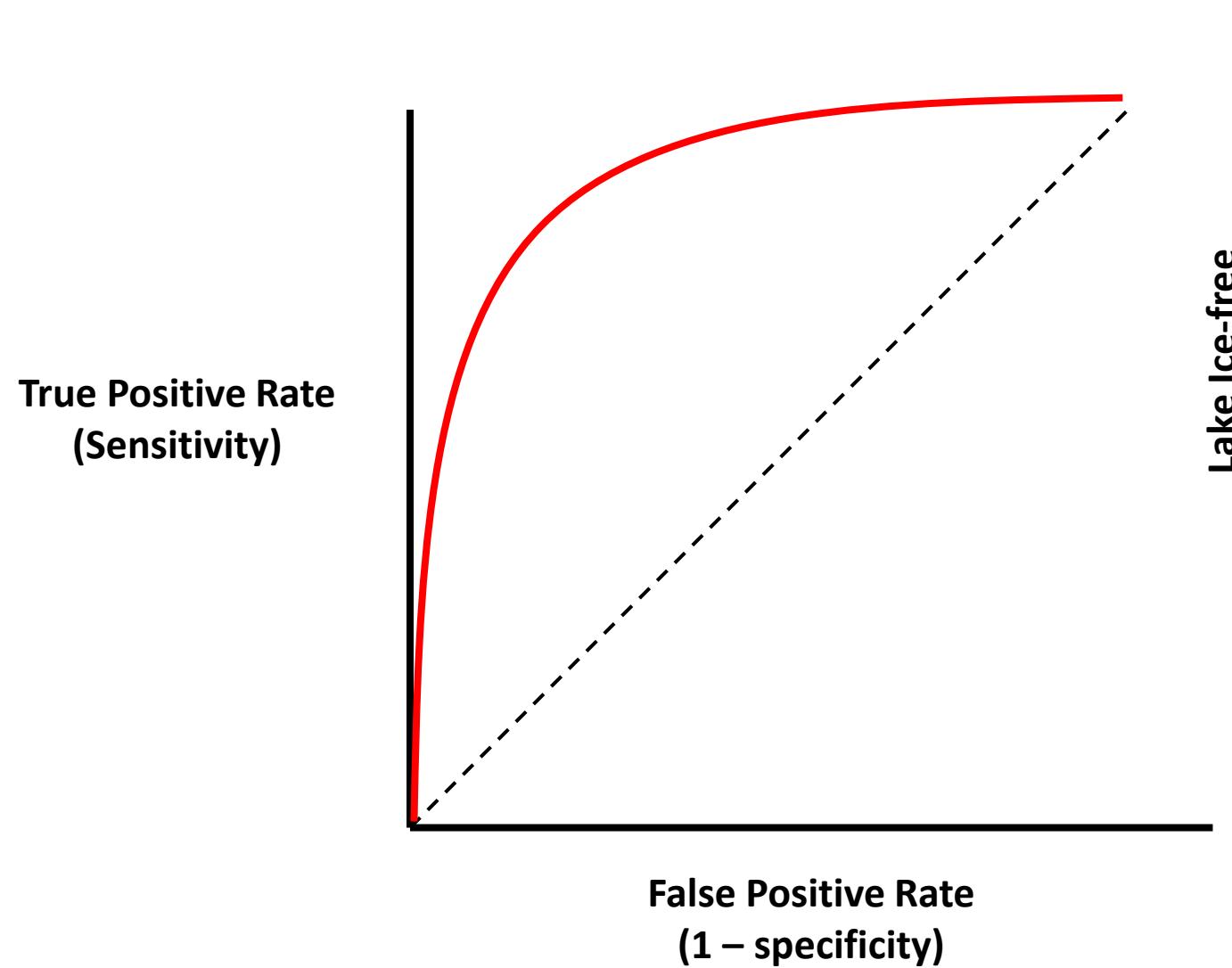
# Optimizing the threshold - ROC



# Optimizing the threshold - ROC



# Optimizing the threshold - ROC



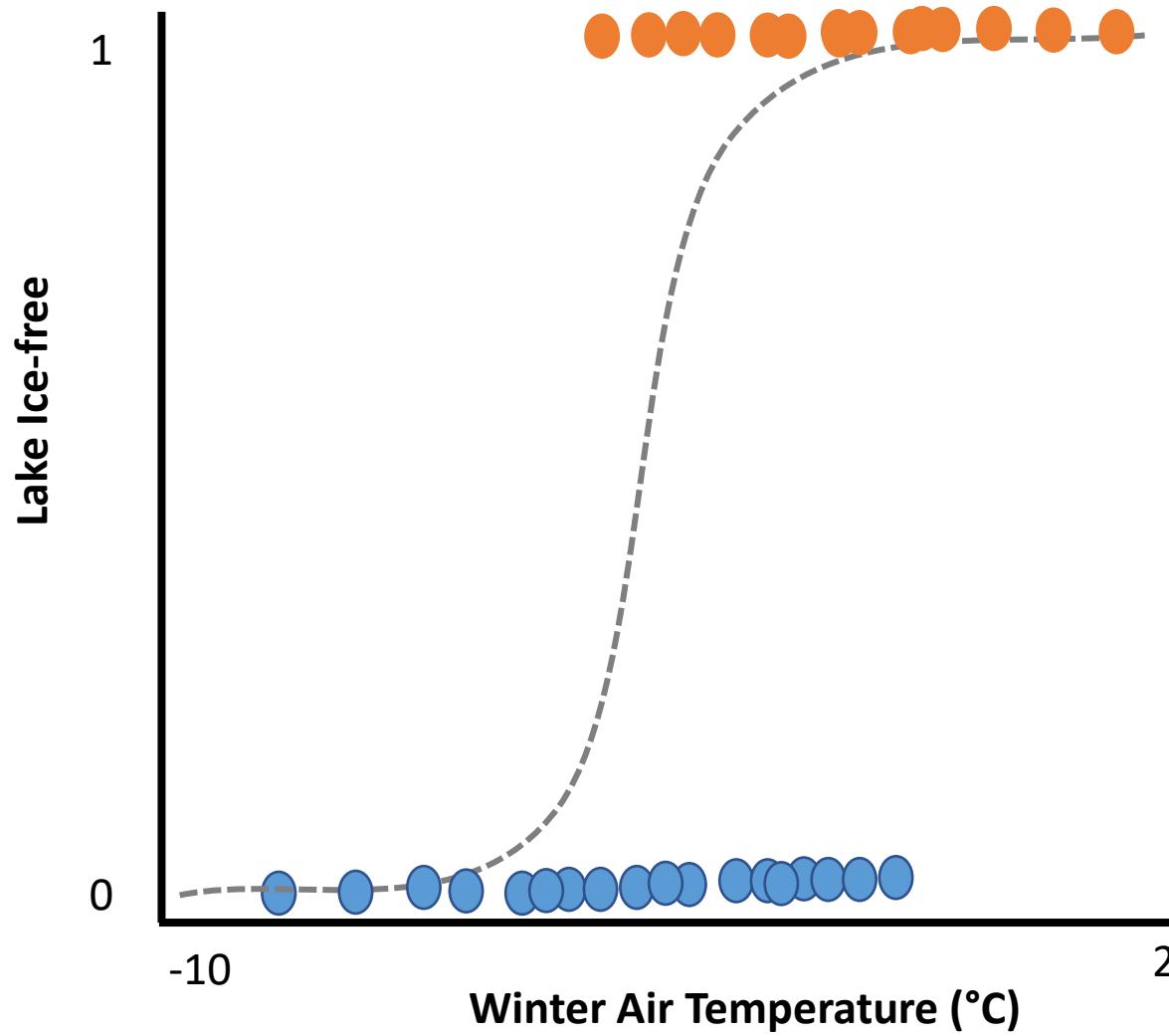
# Comparing models - AUC



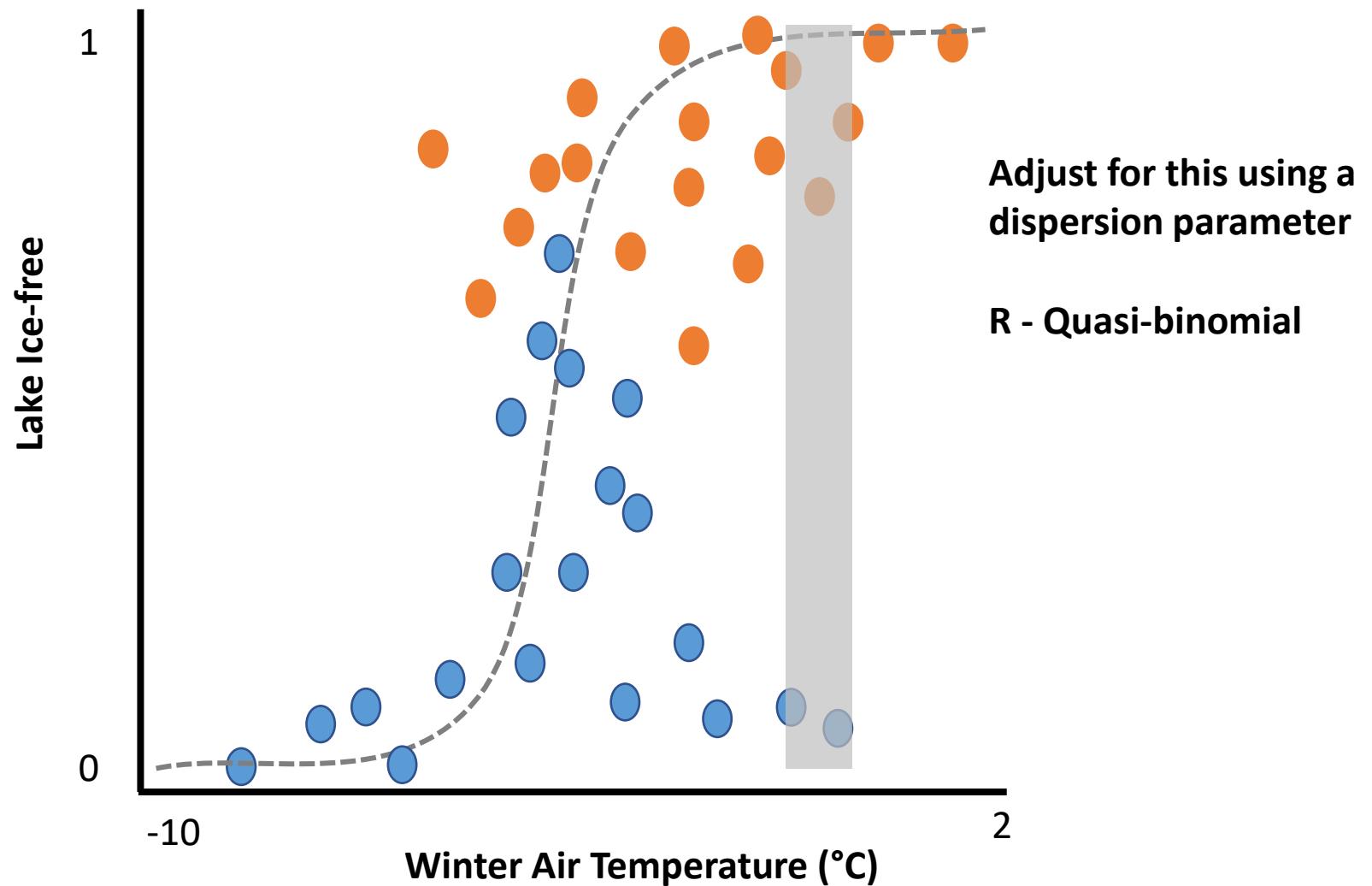
# A note about overdispersion

- Over dispersion when the variance exceeds the mean model fit
- Frequent in ecology data

# A note about overdispersion



# A note about overdispersion



# How do logistic regressions relate to GLMs?

- Use a link function to connect to a linear function
- Use Maximum Likelihood rather than sum of squares
- Allows for the flexibility in GLMs for many distributions

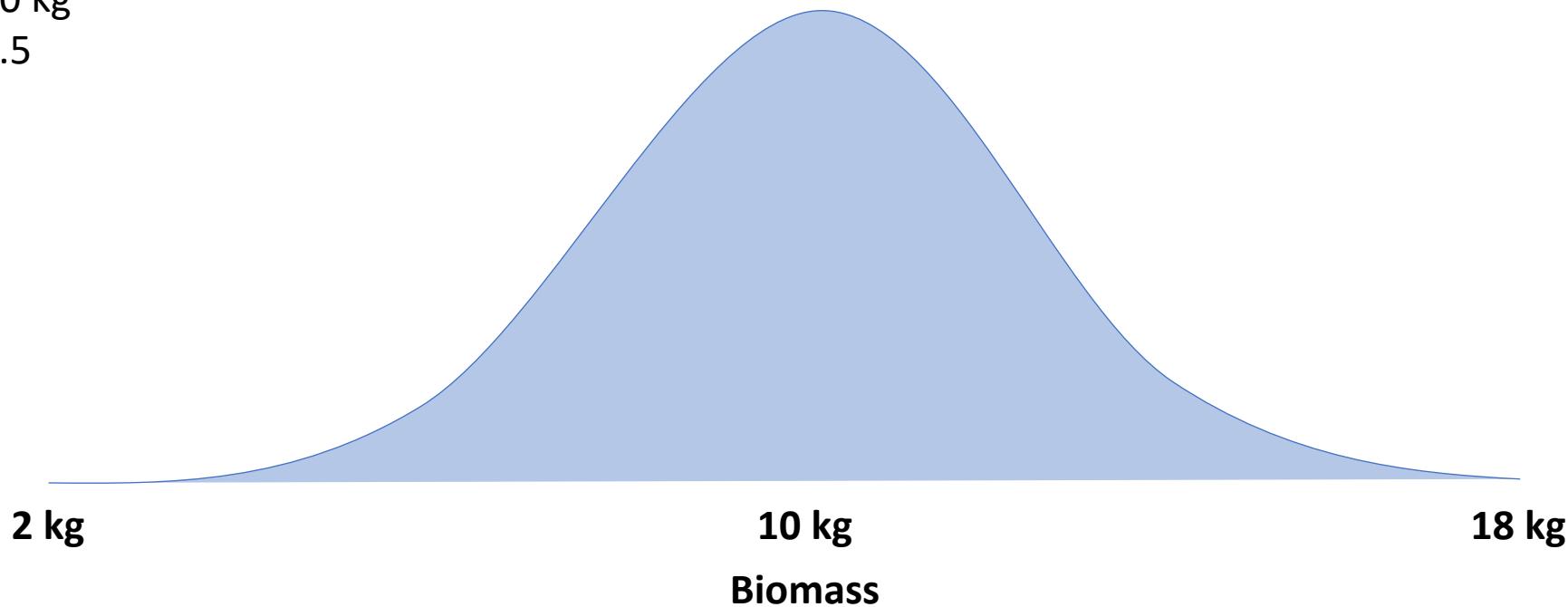
# GLM using a normal distribution

## Probability vs. Likelihood

Normal distribution

$$\mu = 10 \text{ kg}$$

$$\sigma = 3.5$$



# GLM using a normal distribution

## Probability vs. Likelihood

Normal distribution

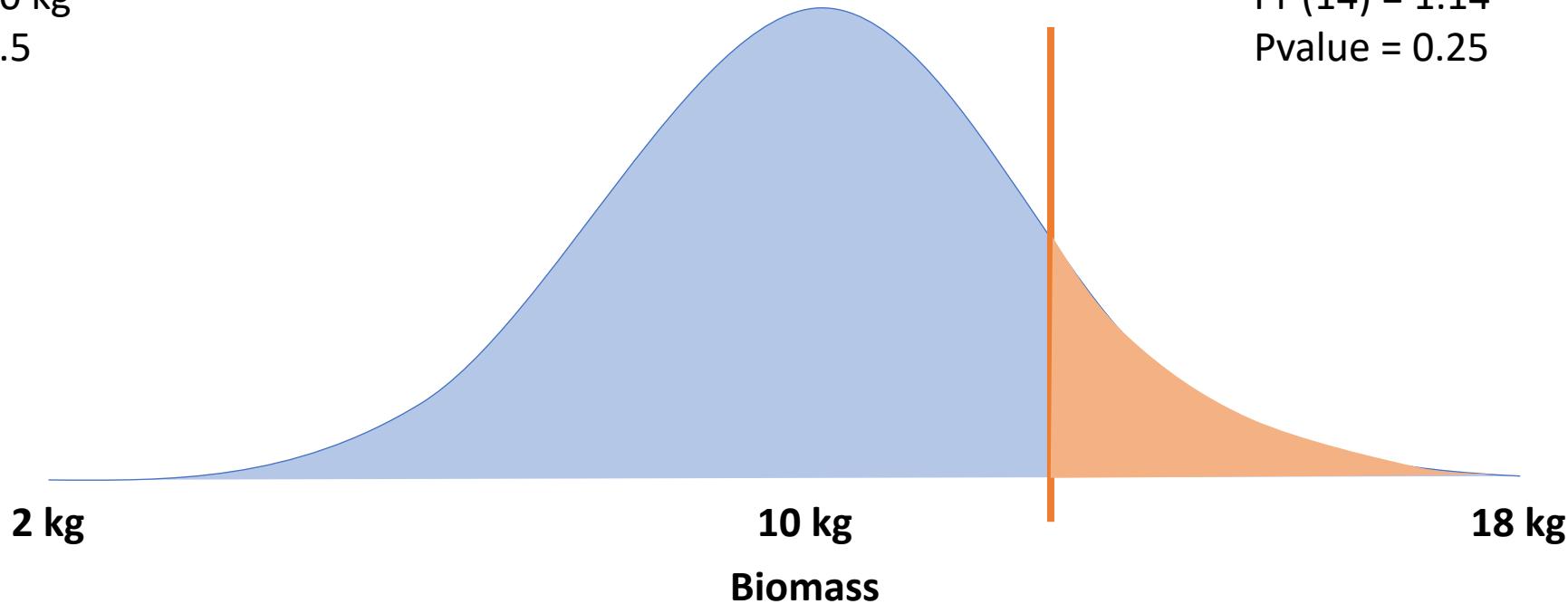
$$\mu = 10 \text{ kg}$$

$$\sigma = 3.5$$

$$\Pr(14) = ?$$

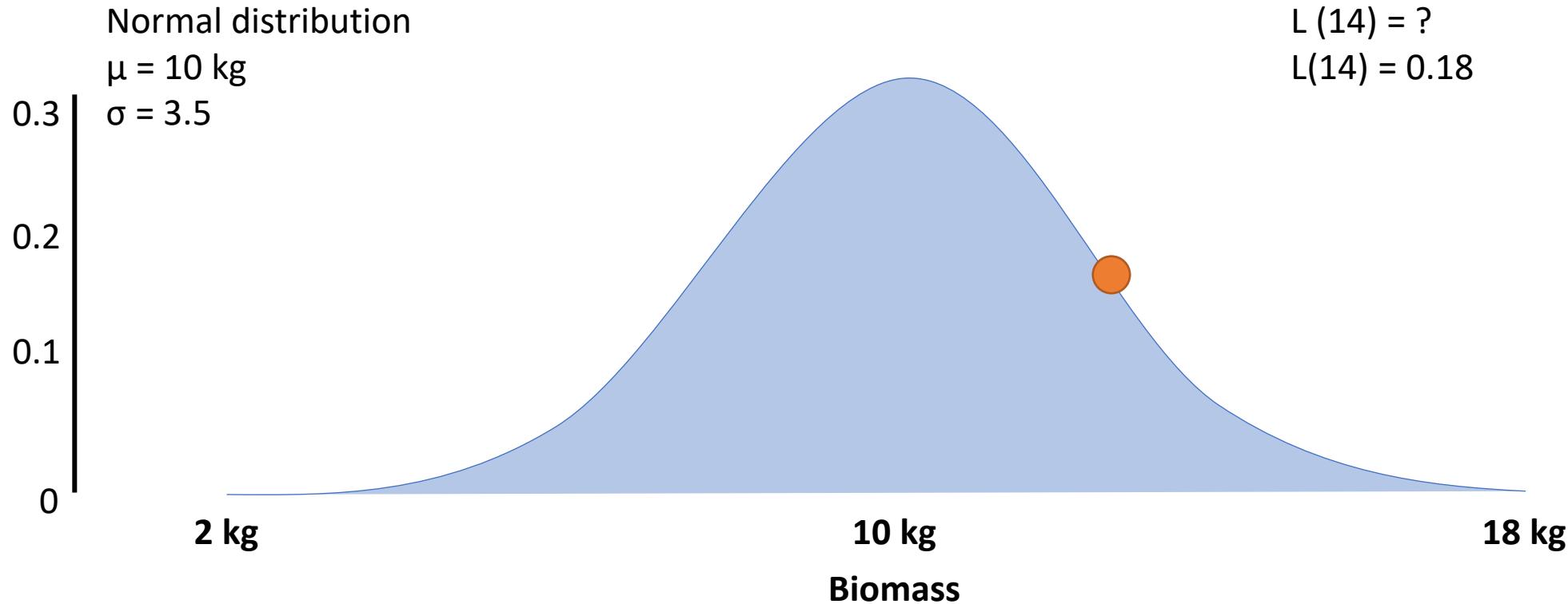
$$\Pr(14) = 1.14$$

$$\text{Pvalue} = 0.25$$



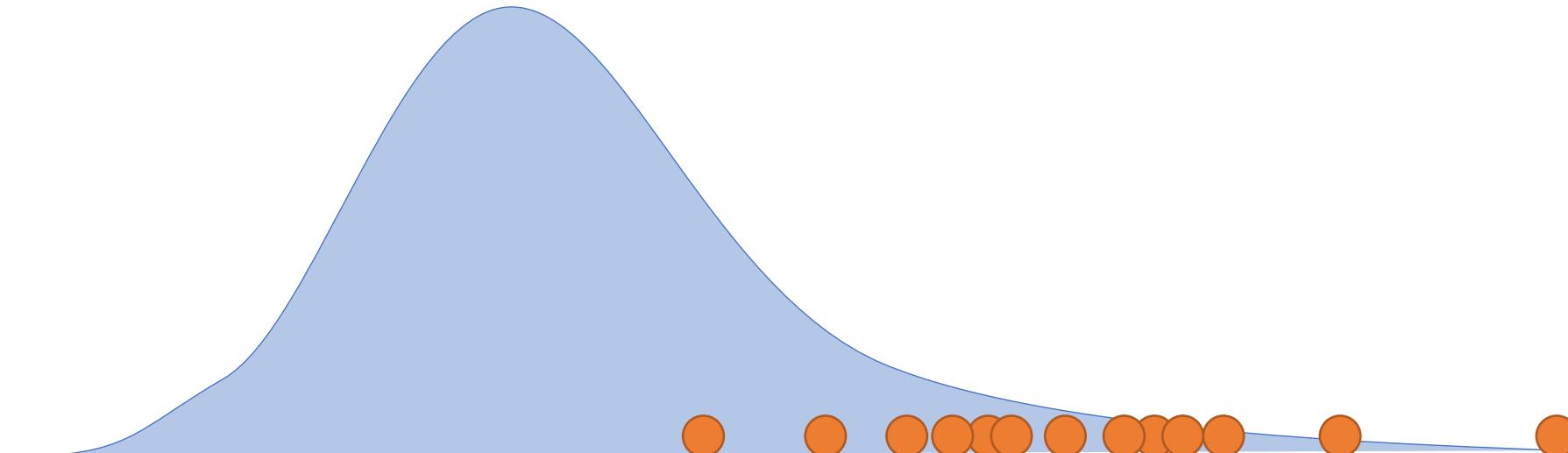
# GLM using a normal distribution

## Probability vs. Likelihood



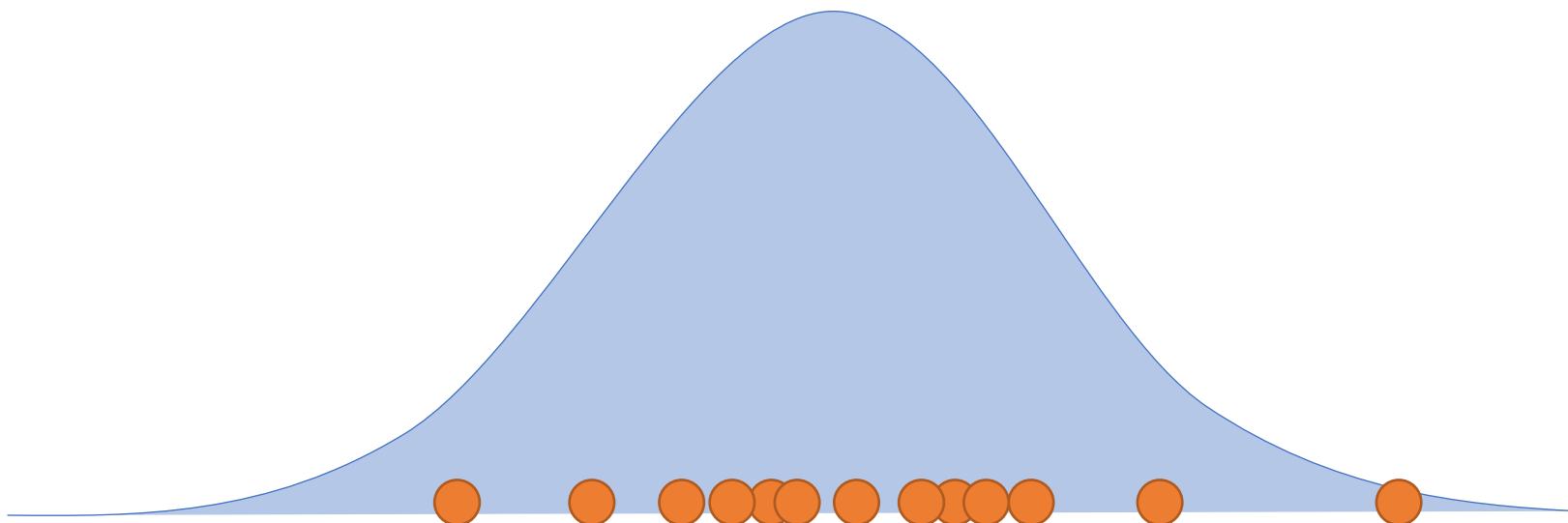
# Maximum likelihood for normal distribution

Maximum likelihood =  $L(1) * L(2) * \dots * L(i)$



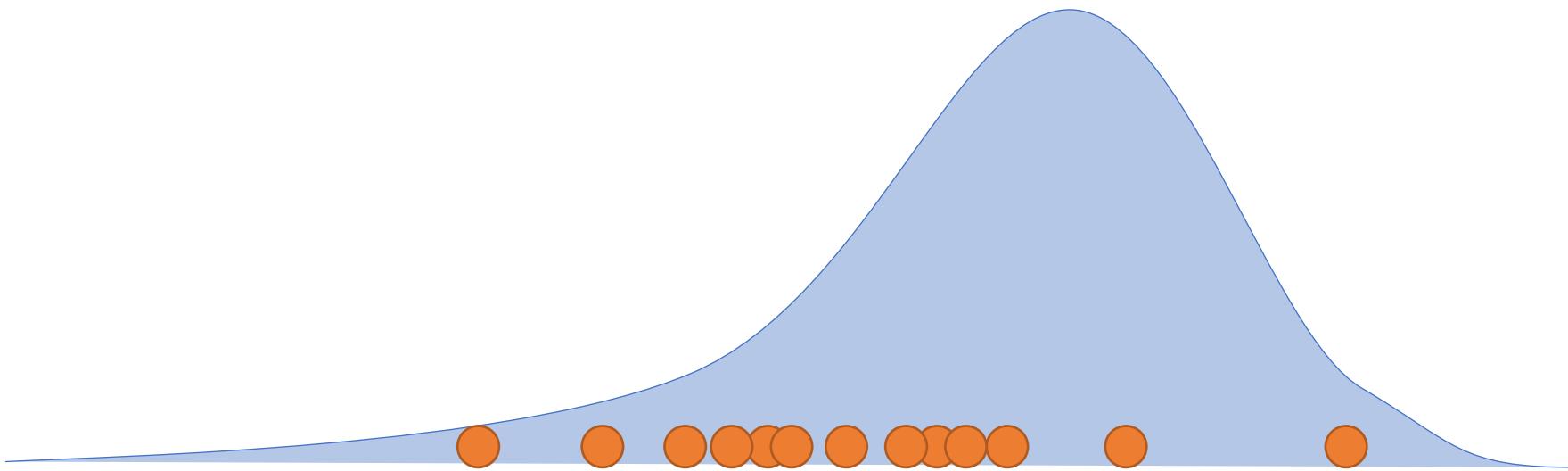
# Maximum likelihood for normal distribution

Maximum likelihood =  $L(1) * L(2) * \dots * L(i)$

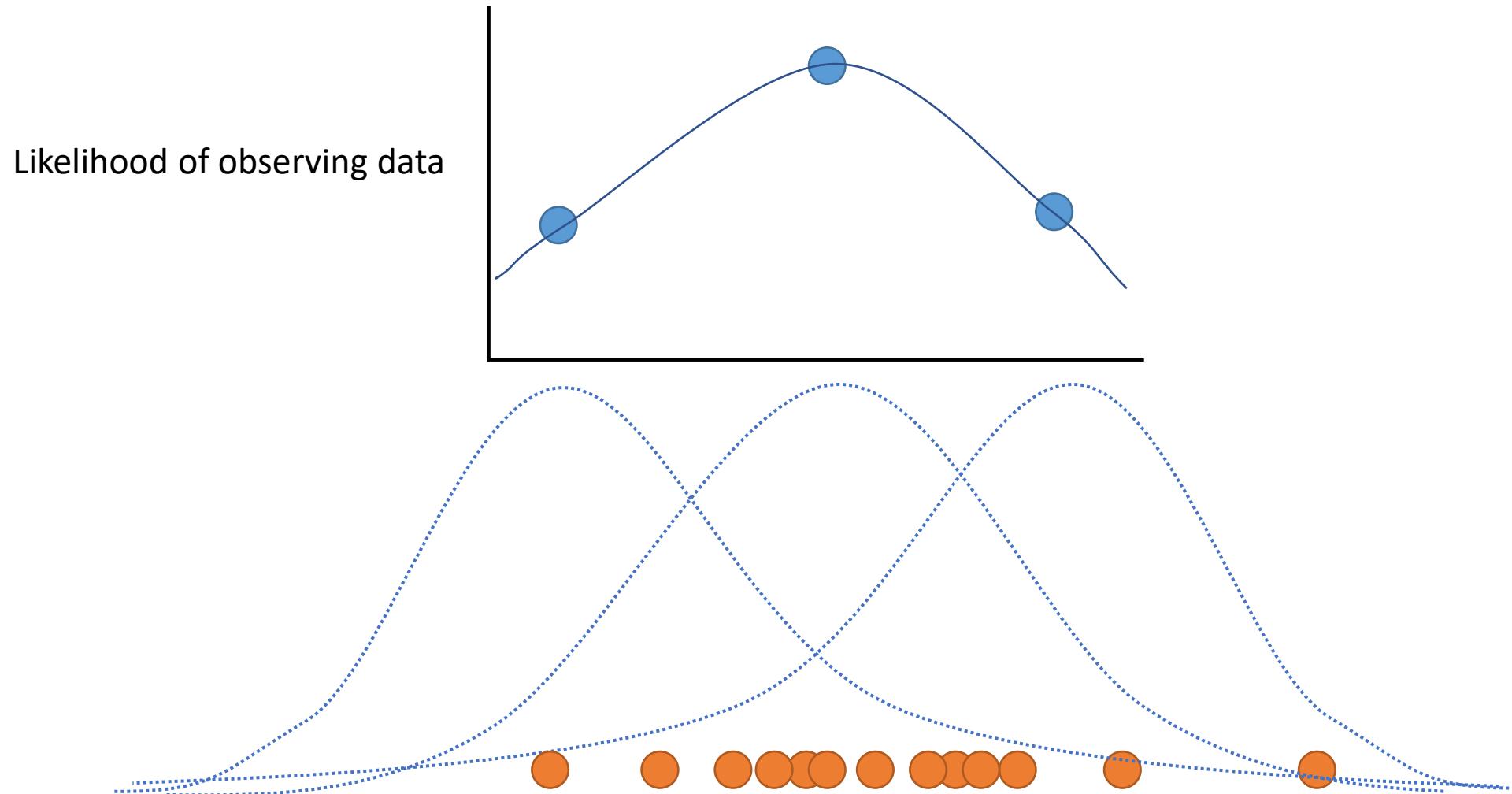


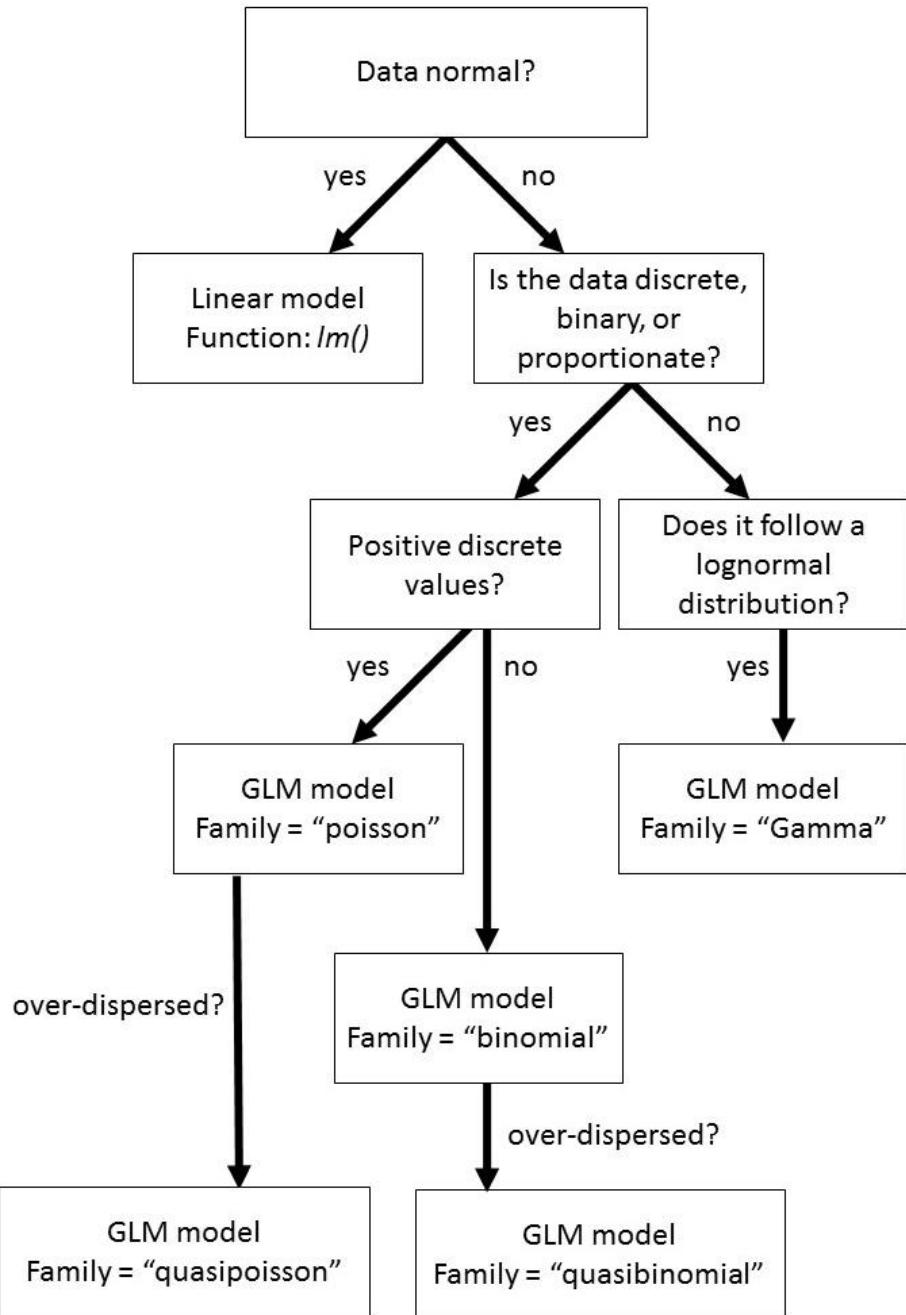
# Maximum likelihood for normal distribution

Maximum likelihood =  $L(1) * L(2) * \dots * L(i)$



# Maximum likelihood for normal distribution



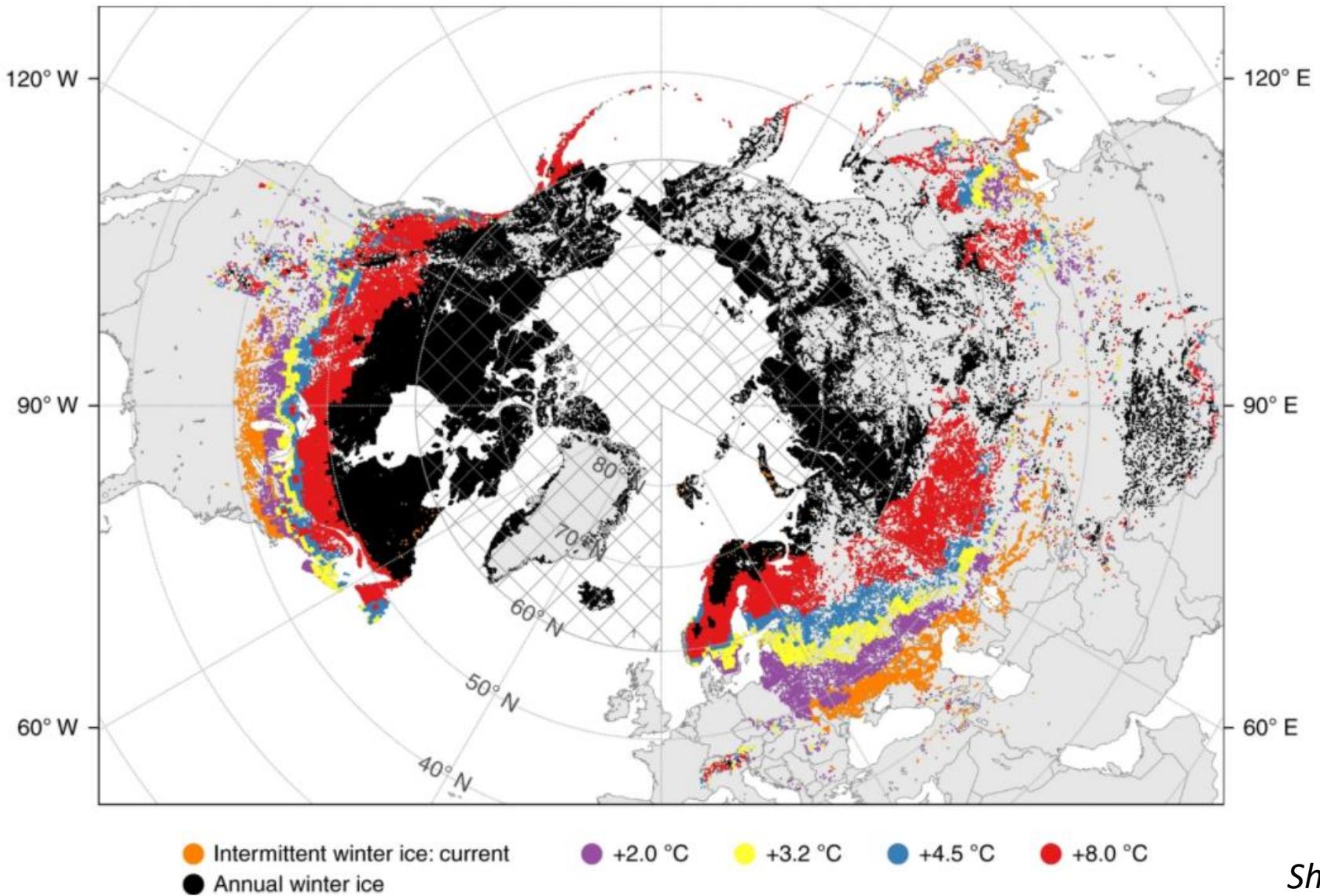


# A simple workflow for GLMs

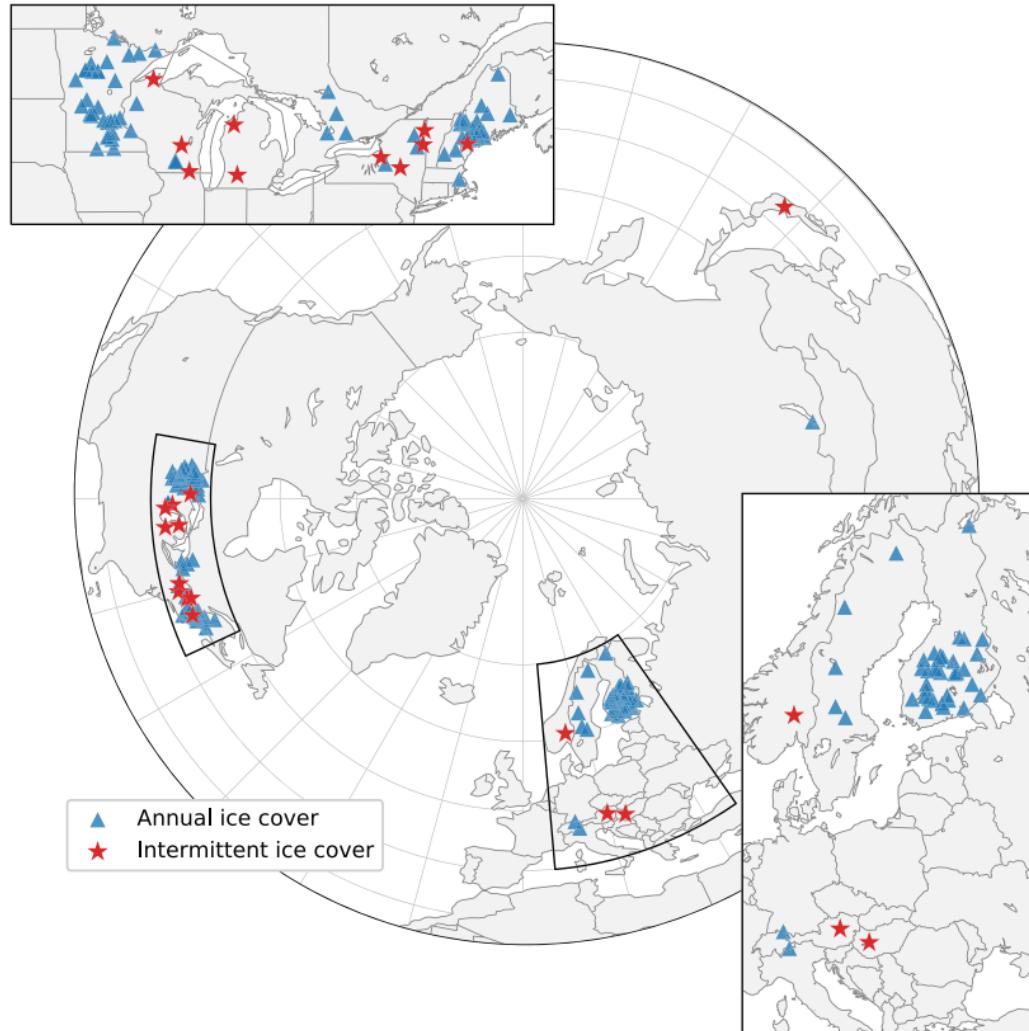


Case Study: Logistic regression to identify extreme events

# Lake ice is threatened by climate change

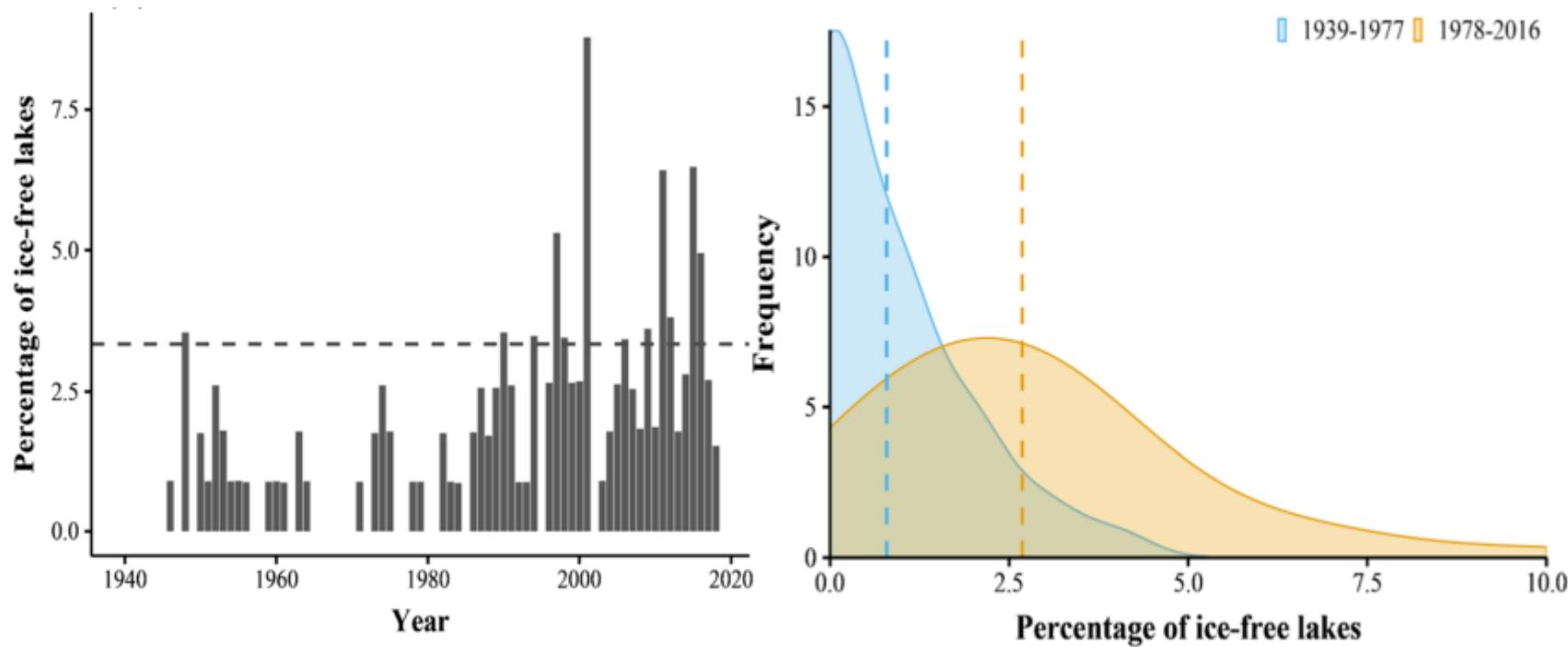


# Extreme events on lake ice



- Selected 122 lakes in the Northern Hemisphere
- Tested the frequency of ice-free years over time
- Examined the role of extreme temperatures

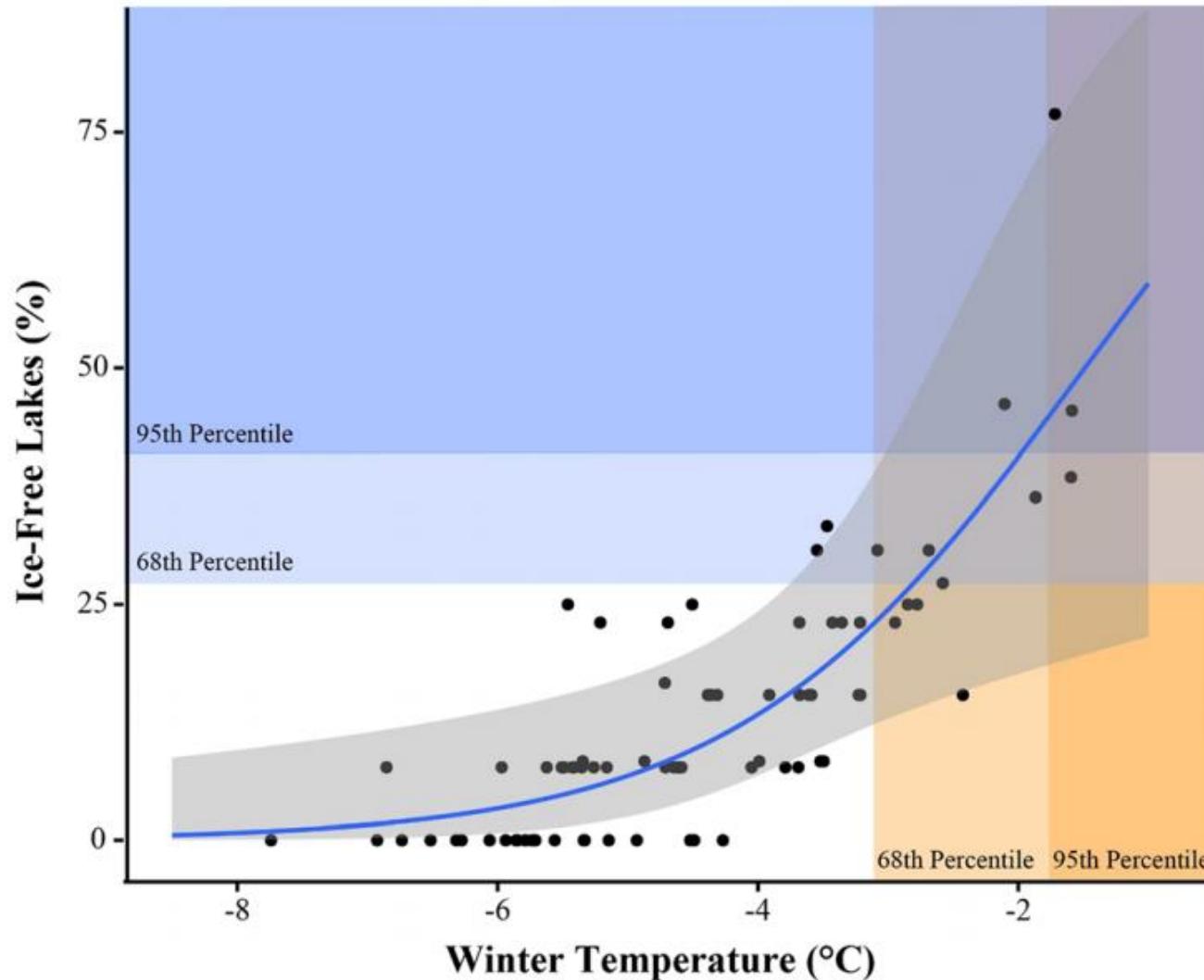
# A large increase in ice-free coverage



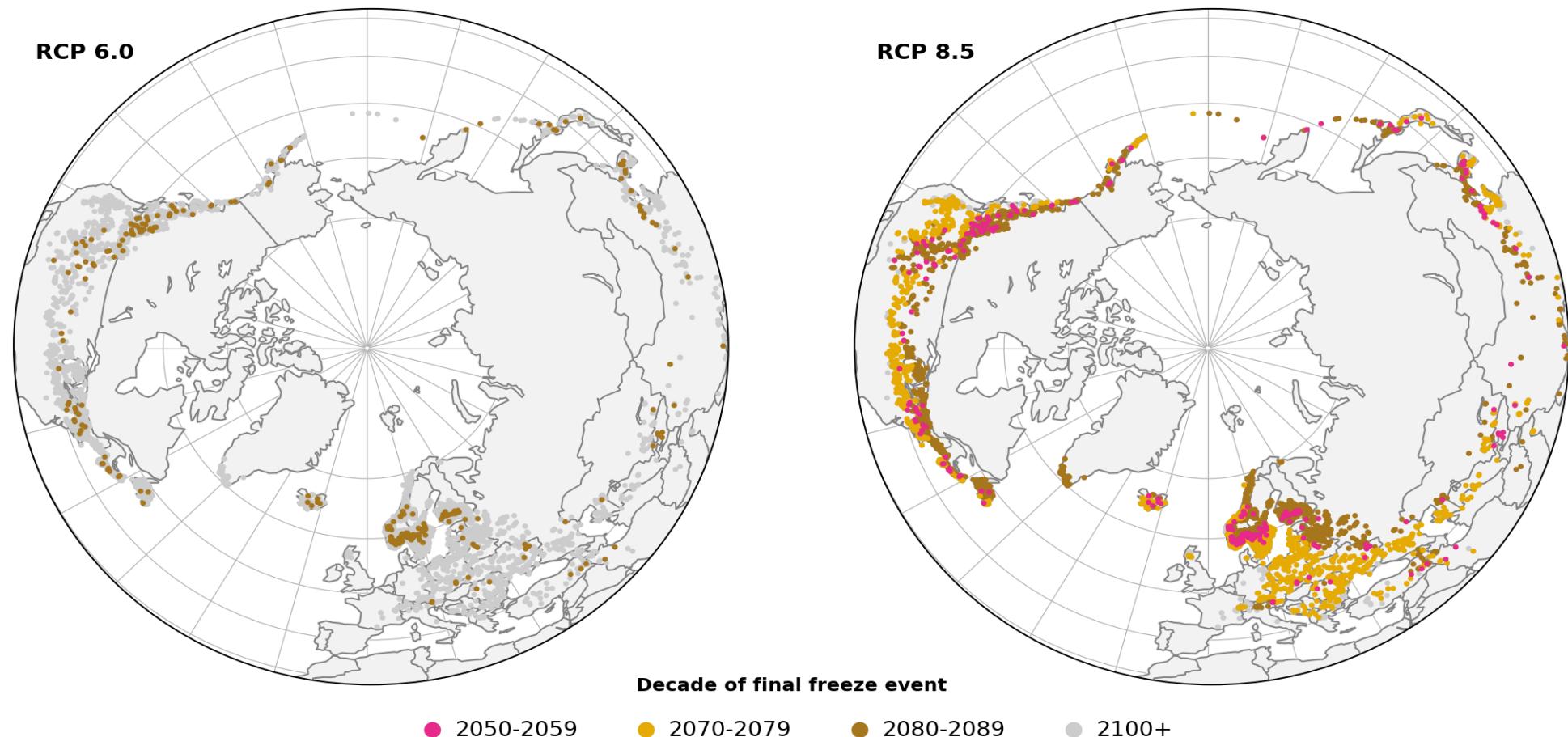
# Data of ice-free events

Lake	Year	IceFree	Winter air temperature
Sebago Lake	1995	1	-1
Sebago Lake	1996	0	-2.5
Sebago Lake	1997	0	-2.9
Sebago Lake	1998	0	-3.1
George Lake	1995	0	-2.1
George Lake	1996	1	-1.6
George Lake	1997	1	-0.5

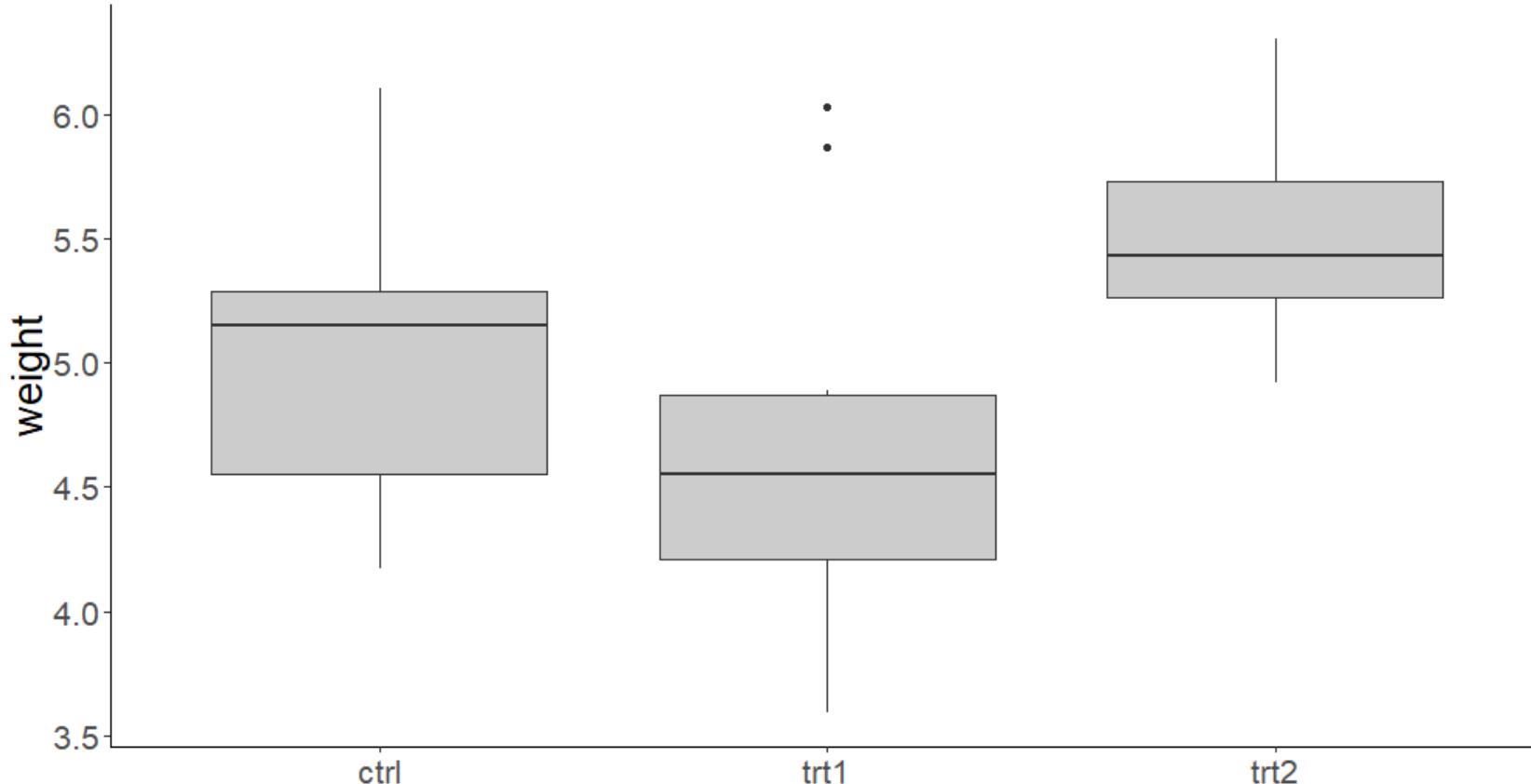
# Extremes in temperature cause extremes in ice coverage



# The future loss of ice coverage



# Post-hoc analyses



**GLM output**

$\chi^2 = 3.77$ , df = 2, p = 0.007

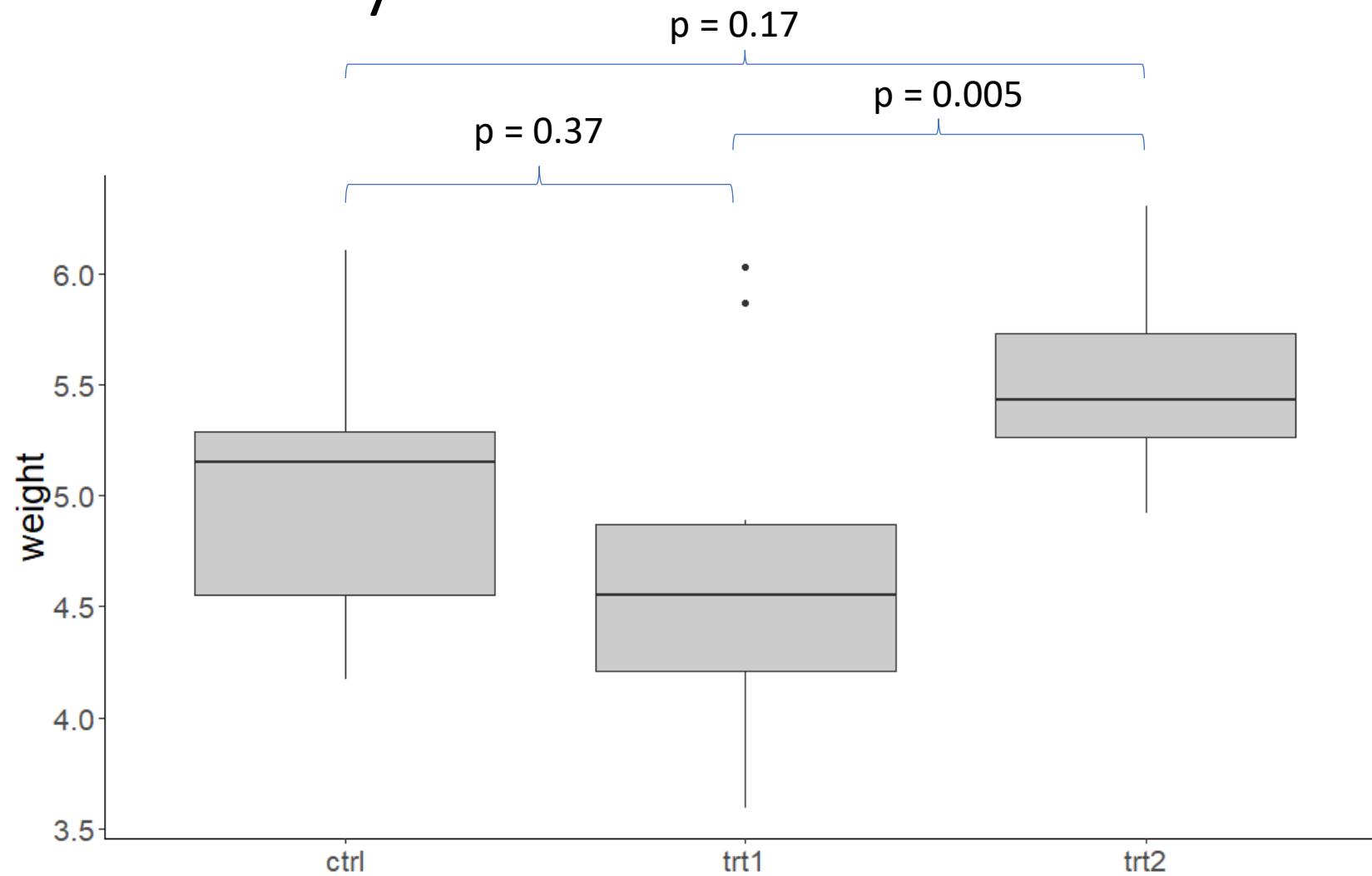
**Which groups are different?**

A – B

B – C

A – C

# Post-hoc analyses



# Generalized linear mixed models

- A GLM that includes both fixed effects and random effects
- **Fixed effects:** predictor to test for significance; constant among observations; estimated using ML in GLMs
- **Random effects:** predictor to account for variation among groups; estimates vary between observations; typically estimated based on a near-normal distribution
- All linear models have fixed effects, but not all have random effects.

# When to use mixed models

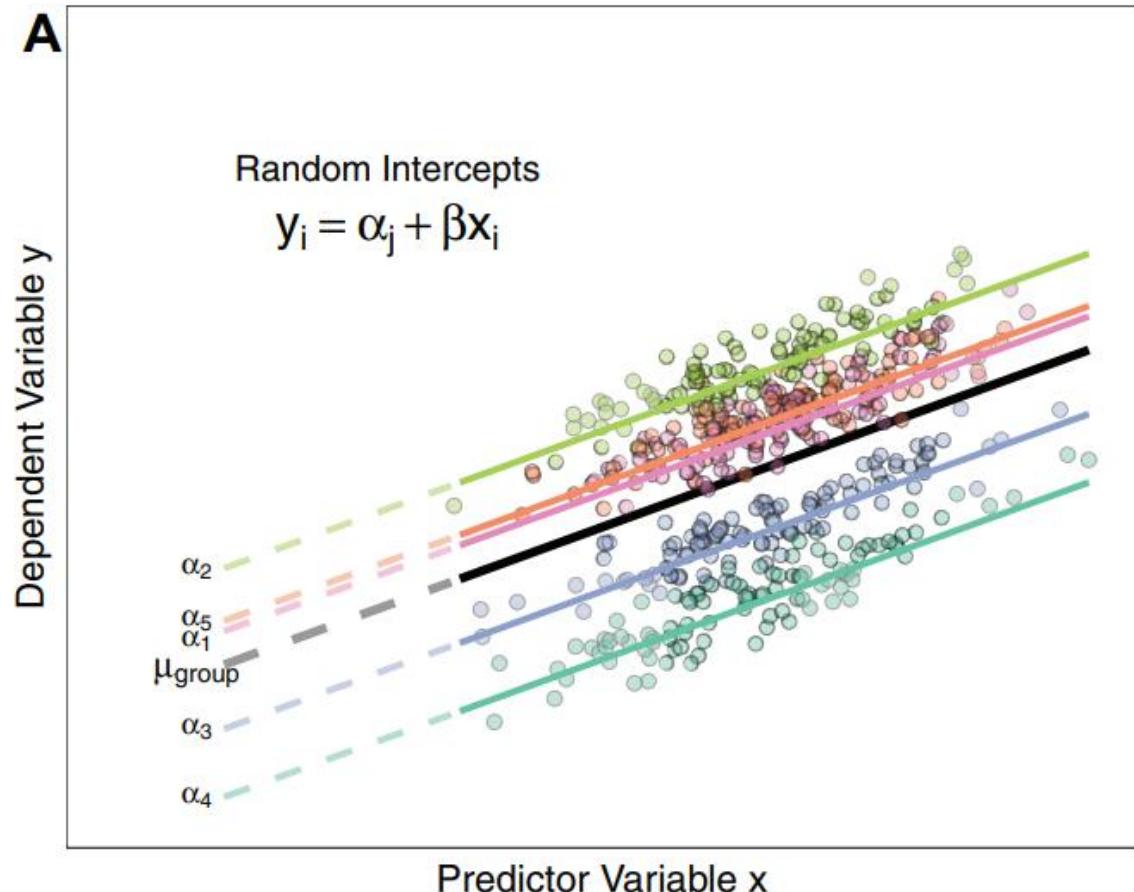
When trying to test for a particular effect while controlling for another variable

Examples:

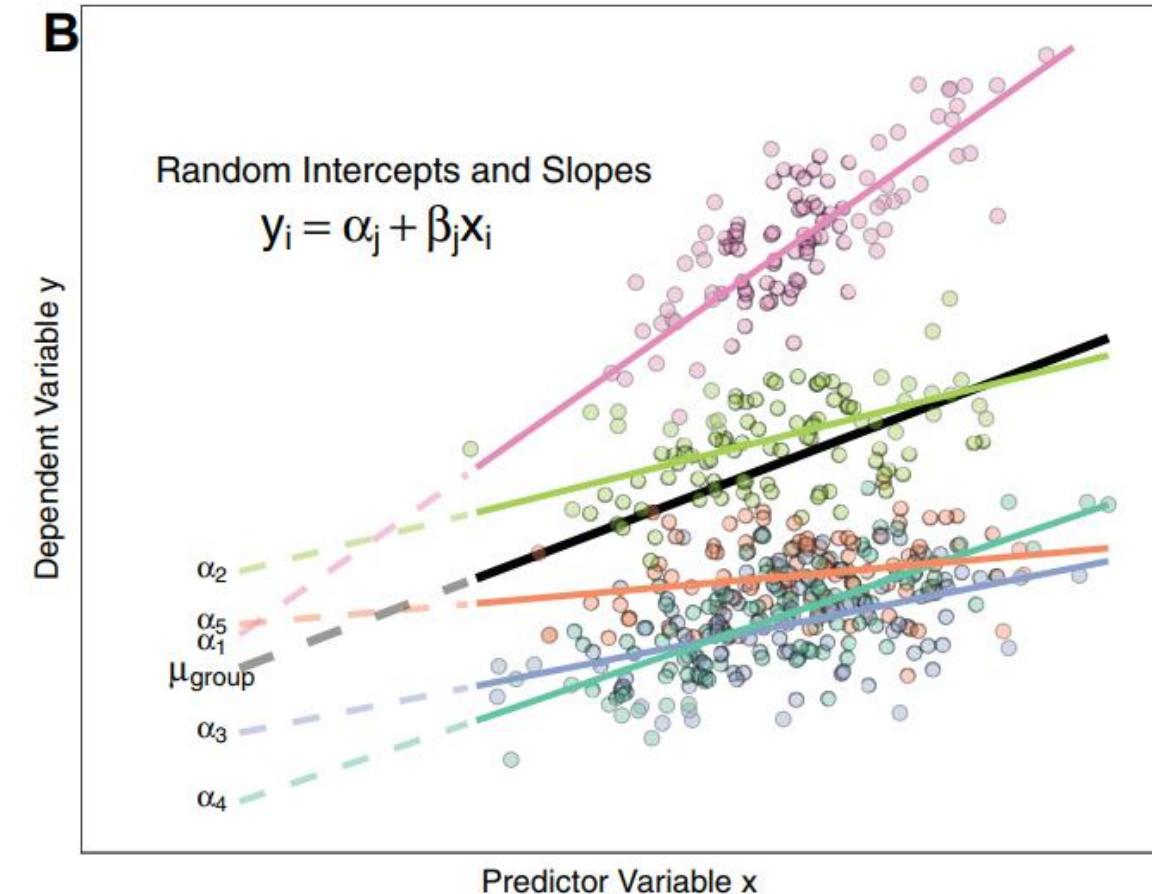
- Testing the effect of a blood pressure drug on patients with different initial BP values
- Testing the effect of nutrient addition on plants surveyed in years with different growing conditions

# Types of mixed models

A



B



# A great introduction to mixed models



## A brief introduction to mixed effects modelling and multi-model inference in ecology

Xavier A. Harrison<sup>1</sup>, Lynda Donaldson<sup>2,3</sup>, Maria Eugenia Correa-Cano<sup>2</sup>, Julian Evans<sup>4,5</sup>, David N. Fisher<sup>4,6</sup>, Cecily E.D. Goodwin<sup>2</sup>, Beth S. Robinson<sup>2,7</sup>, David J. Hodgson<sup>4</sup> and Richard Inger<sup>2,4</sup>

<sup>1</sup> Institute of Zoology, Zoological Society of London, London, UK

<sup>2</sup> Environment and Sustainability Institute, University of Exeter, Penryn, UK

<sup>3</sup> Wildfowl and Wetlands Trust, Slimbridge, Gloucestershire, UK

<sup>4</sup> Centre for Ecology and Conservation, University of Exeter, Penryn, UK

<sup>5</sup> Department of Biology, University of Ottawa, Ottawa, ON, Canada

<sup>6</sup> Department of Integrative Biology, University of Guelph, Guelph, ON, Canada

<sup>7</sup> WildTeam Conservation, Padstow, UK

### ABSTRACT

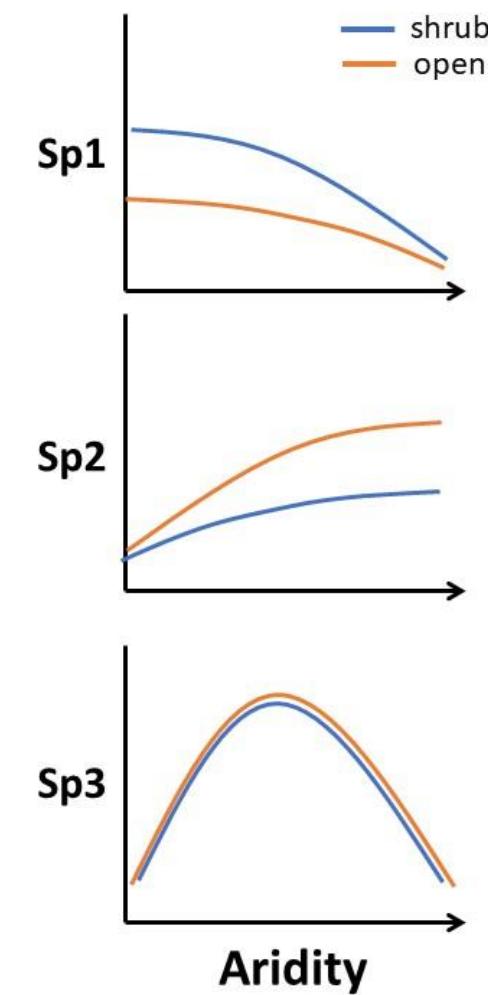
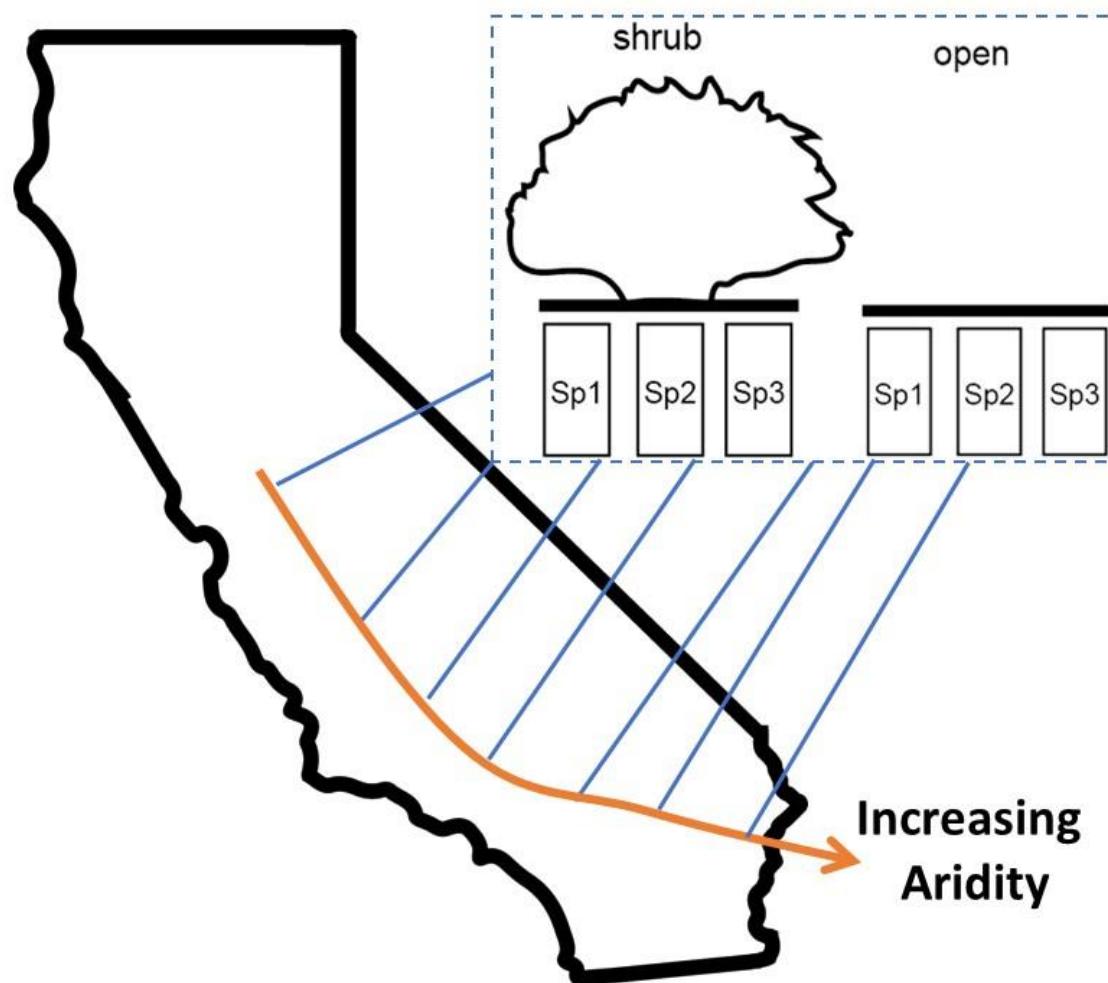
The use of linear mixed effects models (LMMs) is increasingly common in the analysis of biological data. Whilst LMMs offer a flexible approach to modelling a broad range of data types, ecological data are often complex and require complex model structures, and the fitting and interpretation of such models is not always straightforward. The ability to achieve robust biological inference requires that practitioners know how and when to apply these tools. Here, we provide a general overview of current methods for the application of LMMs to biological data, and

<https://peerj.com/articles/4794/>

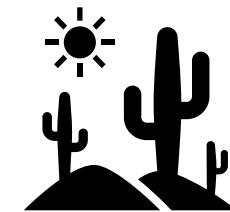


Case Study: Mixed models to estimate plant responses to aridity

# Aridity gradient to test plant interactions



Conducted in 2016



Conducted in 2017

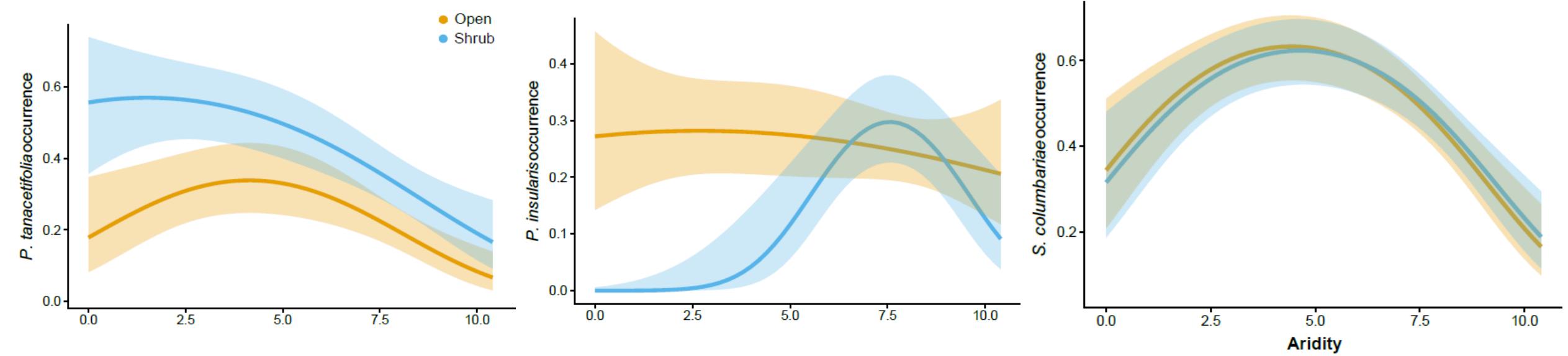


# Data of plant responses to aridity along gradient

PlotID	Year	Site	Aridity	Microsite	Phacelia	Plantago	Salvia
1	2016	Barstow	4.23	open	0	0	1
1	2016	Barstow	4.23	shrub	0	1	1
2	2016	Barstow	4.23	open	1	0	0
2	2016	Barstow	4.23	shrub	0	0	0
3	2016	Barstow	4.23	open	0	0	1
3	2016	Barstow	4.23	shrub	0	0	0
4	2016	Barstow	4.23	open	1	0	1
4	2016	Barstow	4.23	shrub	0	0	0
5	2016	Barstow	4.23	open	0	0	0
5	2016	Barstow	4.23	shrub	0	1	1
6	2016	Barstow	4.23	open	0	1	1

# Model formulation

```
glmer(plantOcc ~ aridity * microsite + (1|Year), data=data, family="binomial")
```







# UNIVERSITY OF TORONTO

Thank you!

<https://www.filazzola.info/>



<https://afilazzola.github.io/IntroGLM/>