# Linear Classification - Logistic Regression
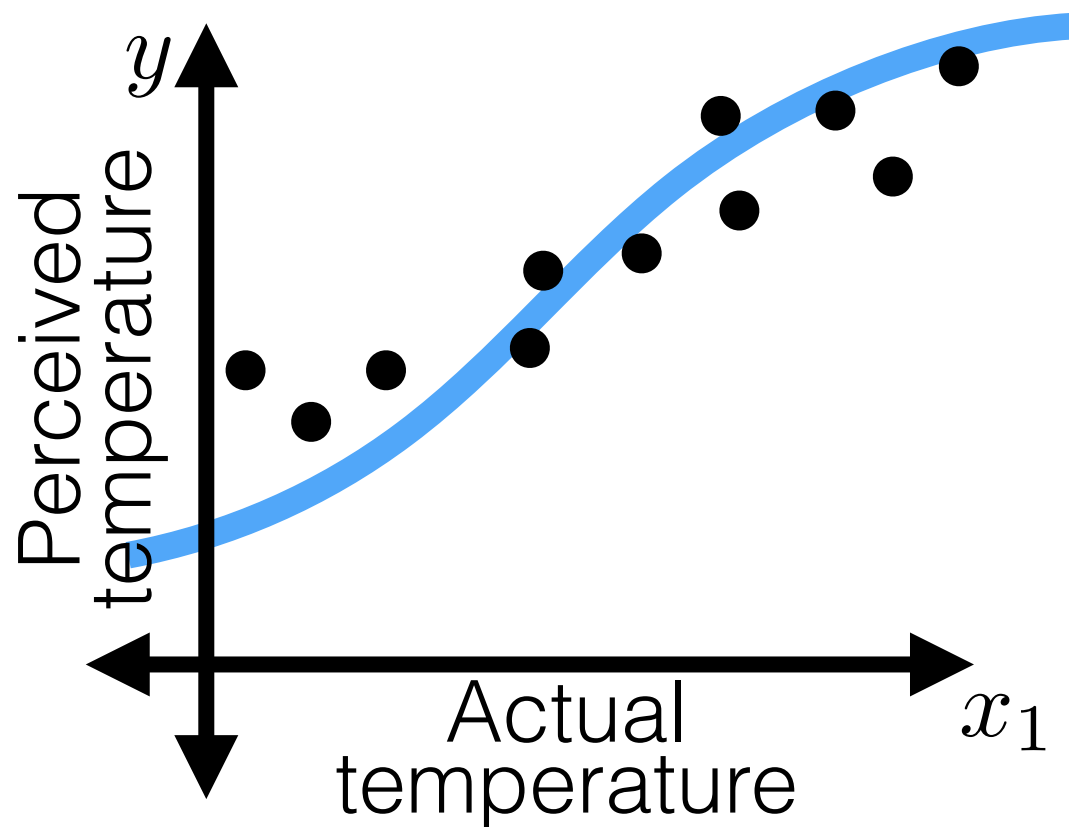
## Prof. Tamara Broderick
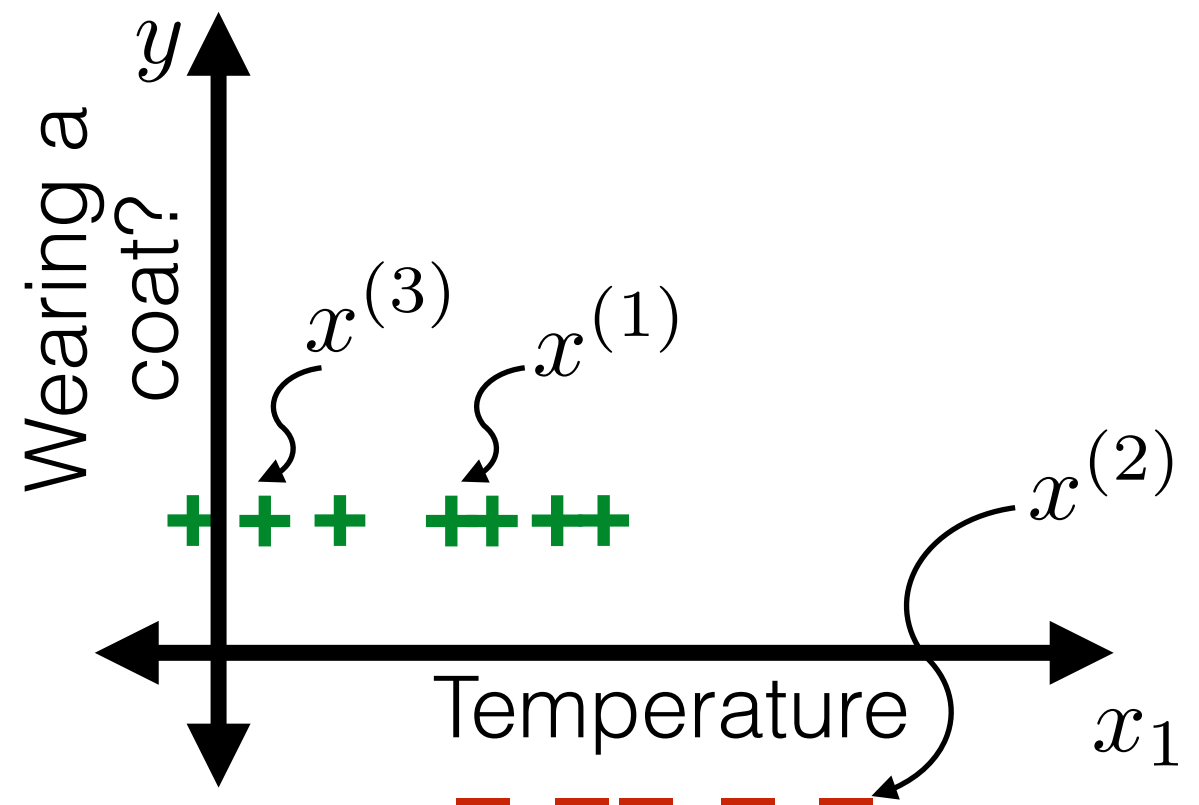
# Recall

## Regression

- Datum $i$: feature vector
  $x^{(i)} = (x_1^{(i)}, \ldots, x_d^{(i)})^\top \in \mathbb{R}^d$

  - Label $y^{(i)} \in \mathbb{R}$

- Hypothesis $h : \mathbb{R}^d \to \mathbb{R}$
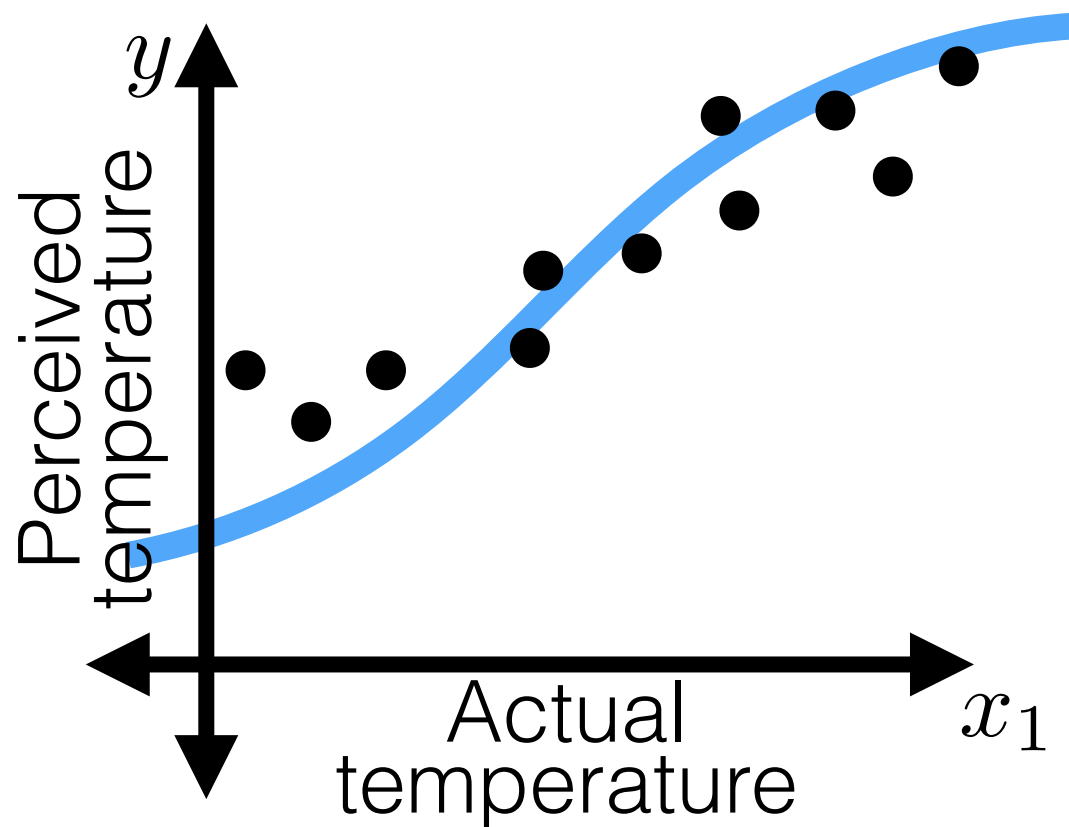


# Compare

## (Two-class) Classification

- Datum $i$: feature vector
  $x^{(i)} = (x_1^{(i)}, \ldots, x_d^{(i)})^\top \in \mathbb{R}^d$

  - Label $y^{(i)} \in \{-1, +1\}$

- Hypothesis $h : \mathbb{R}^d \to \{-1, +1\}$
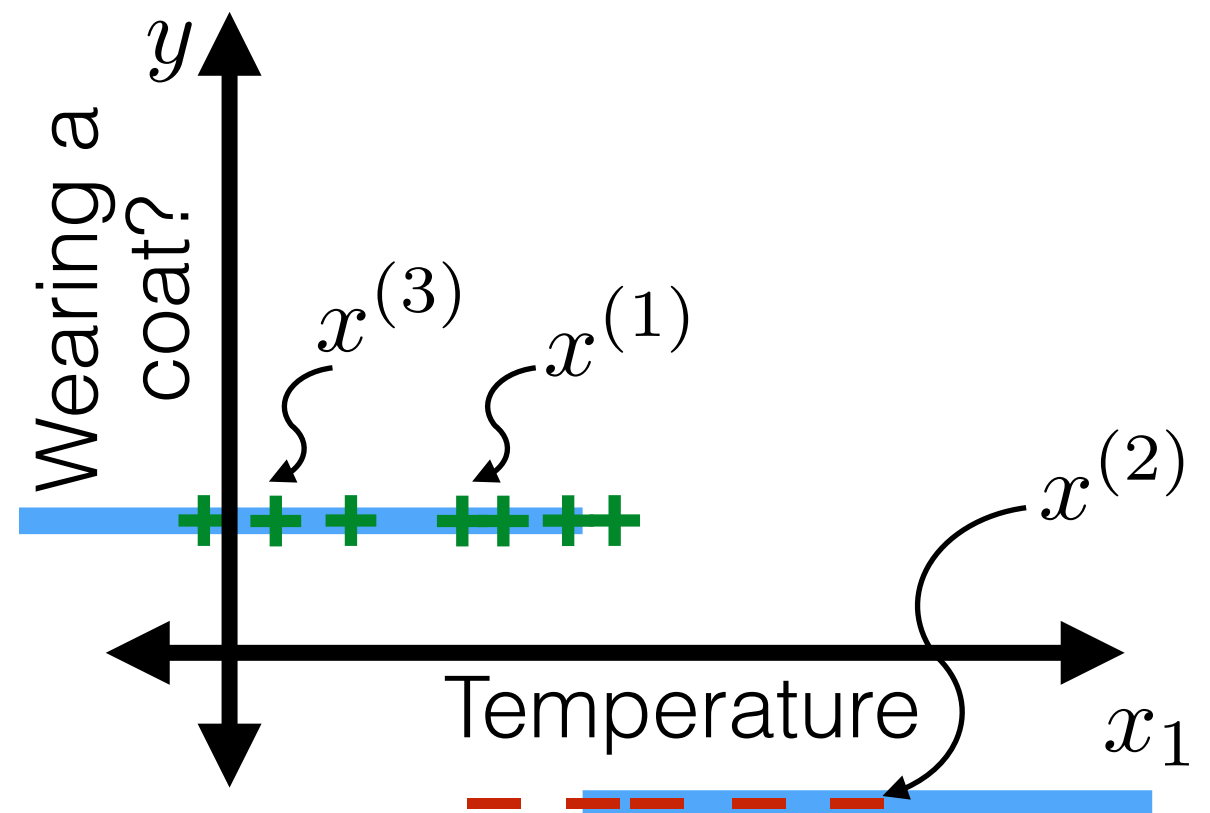
# Recall

## Regression

- Datum $i$: feature vector
$$x^{(i)} = (x_1^{(i)}, \ldots, x_d^{(i)})^\top \in \mathbb{R}^d$$
  - Label $y^{(i)} \in \mathbb{R}$
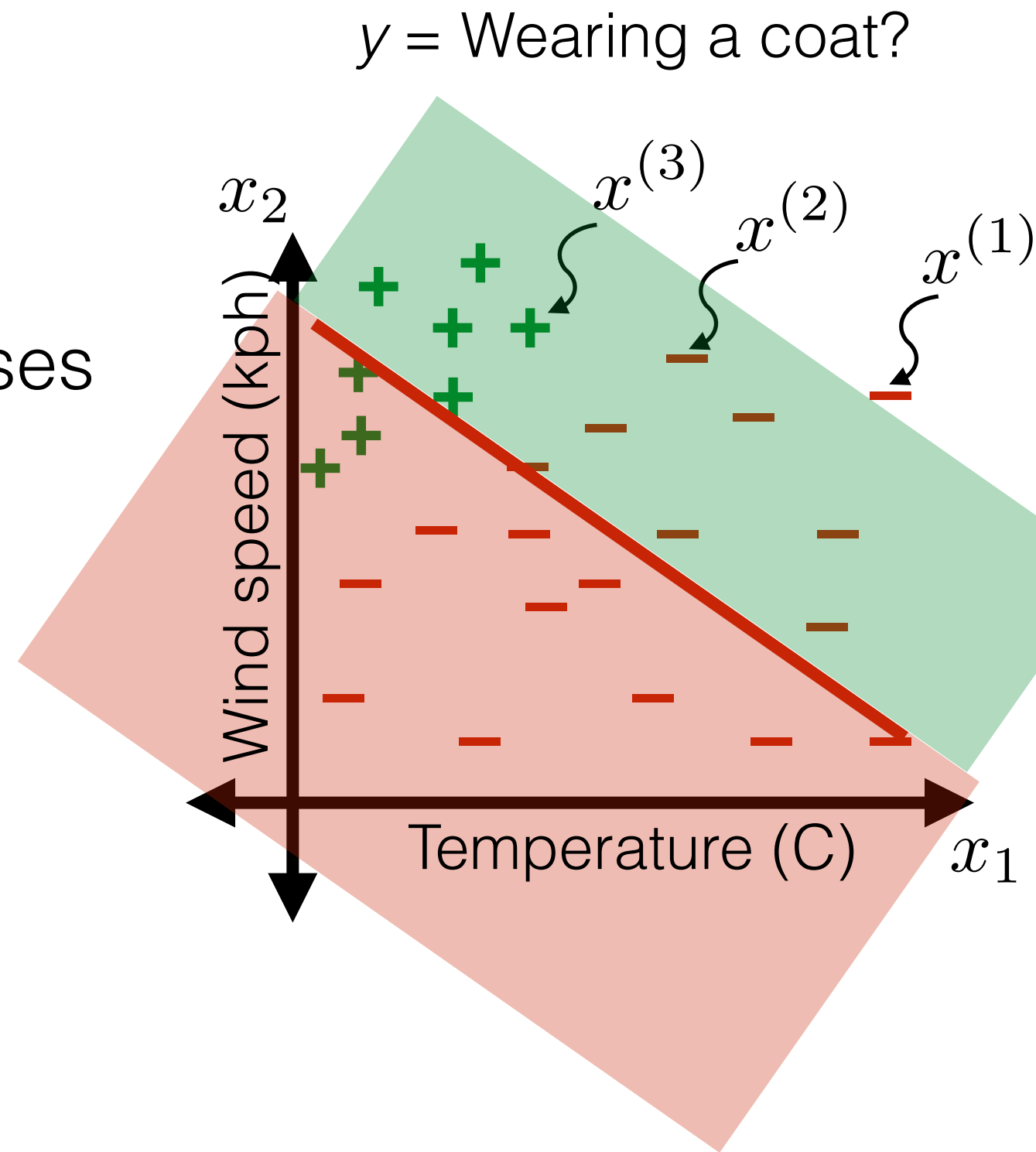- Hypothesis $h : \mathbb{R}^d \to \mathbb{R}$



# Compare

## (Two-class) Classification

- Datum $i$: feature vector
$$x^{(i)} = (x_1^{(i)}, \ldots, x_d^{(i)})^\top \in \mathbb{R}^d$$
  - Label $y^{(i)} \in \{-1, +1\}$
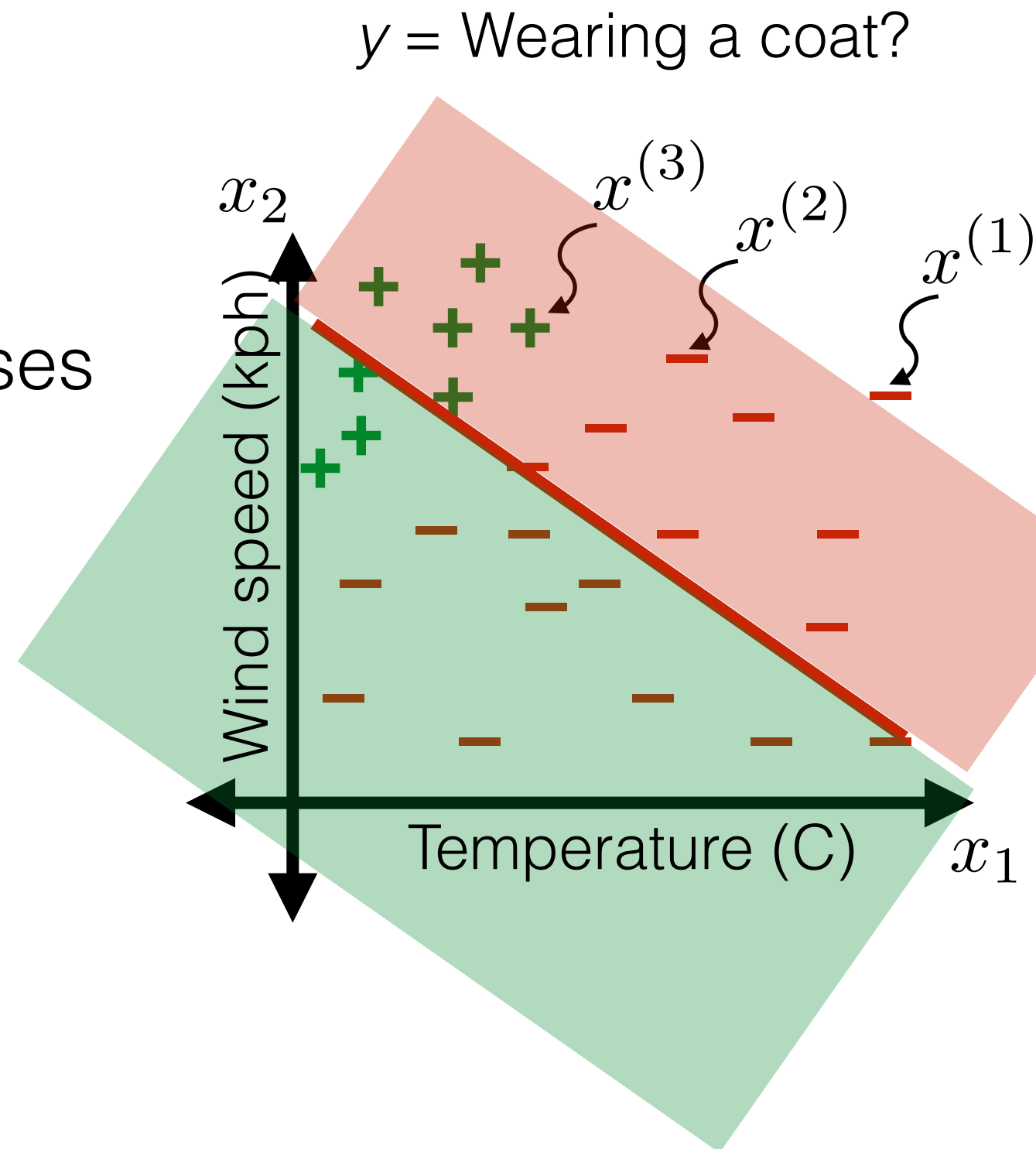- Hypothesis $h : \mathbb{R}^d \to \{-1, +1\}$

# Linear classifiers

- Classification hypothesis:
  $$h : \mathbb{R}^d \to \{-1, +1\}$$
- Linear classifiers $\mathcal{H}$: Hypotheses that label +1 on one side of a line & -1 on the other side

*y* = Wearing a coat?

$x^{(3)}$   $x^{(2)}$   $x^{(1)}$

$x_2$

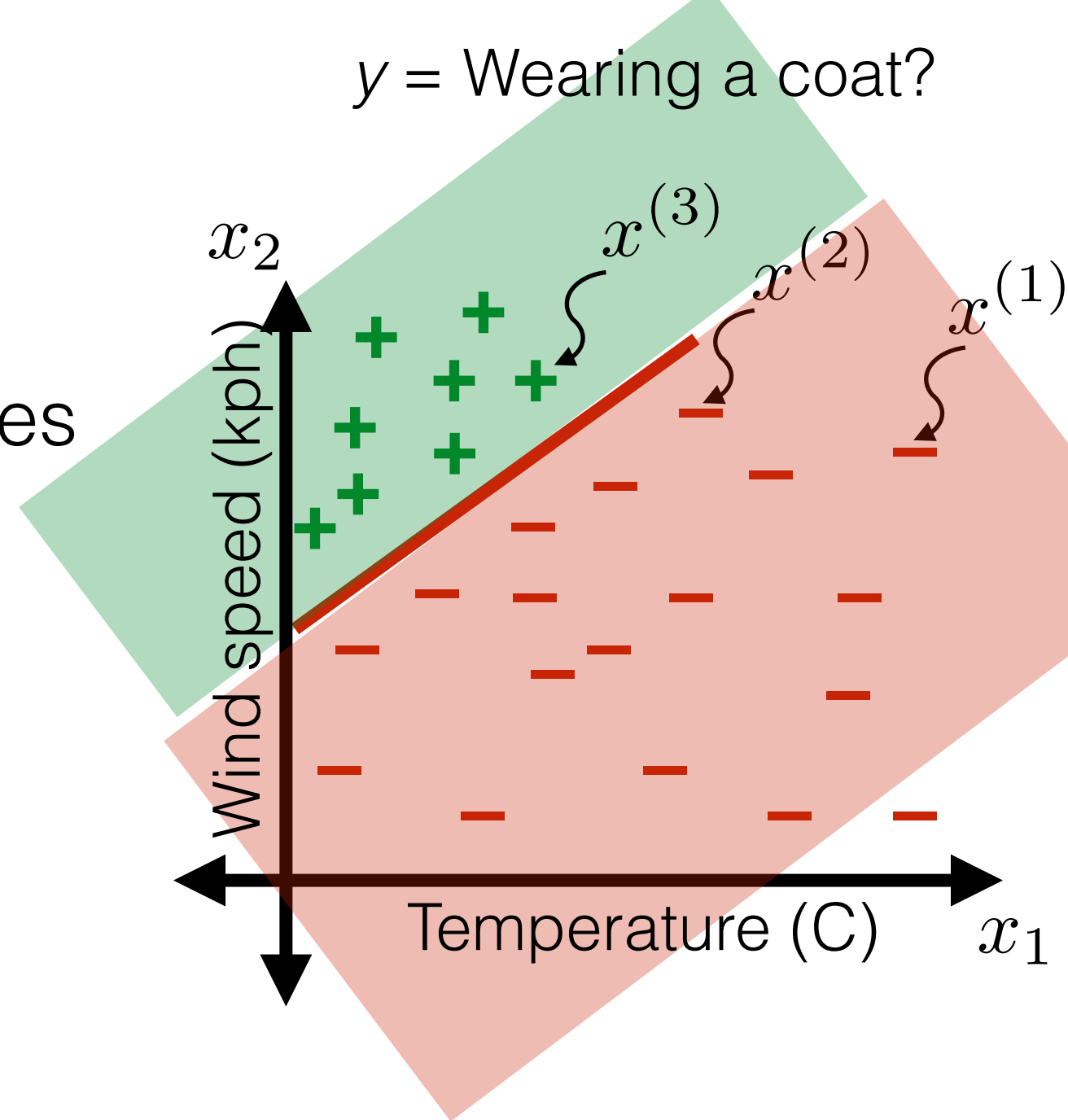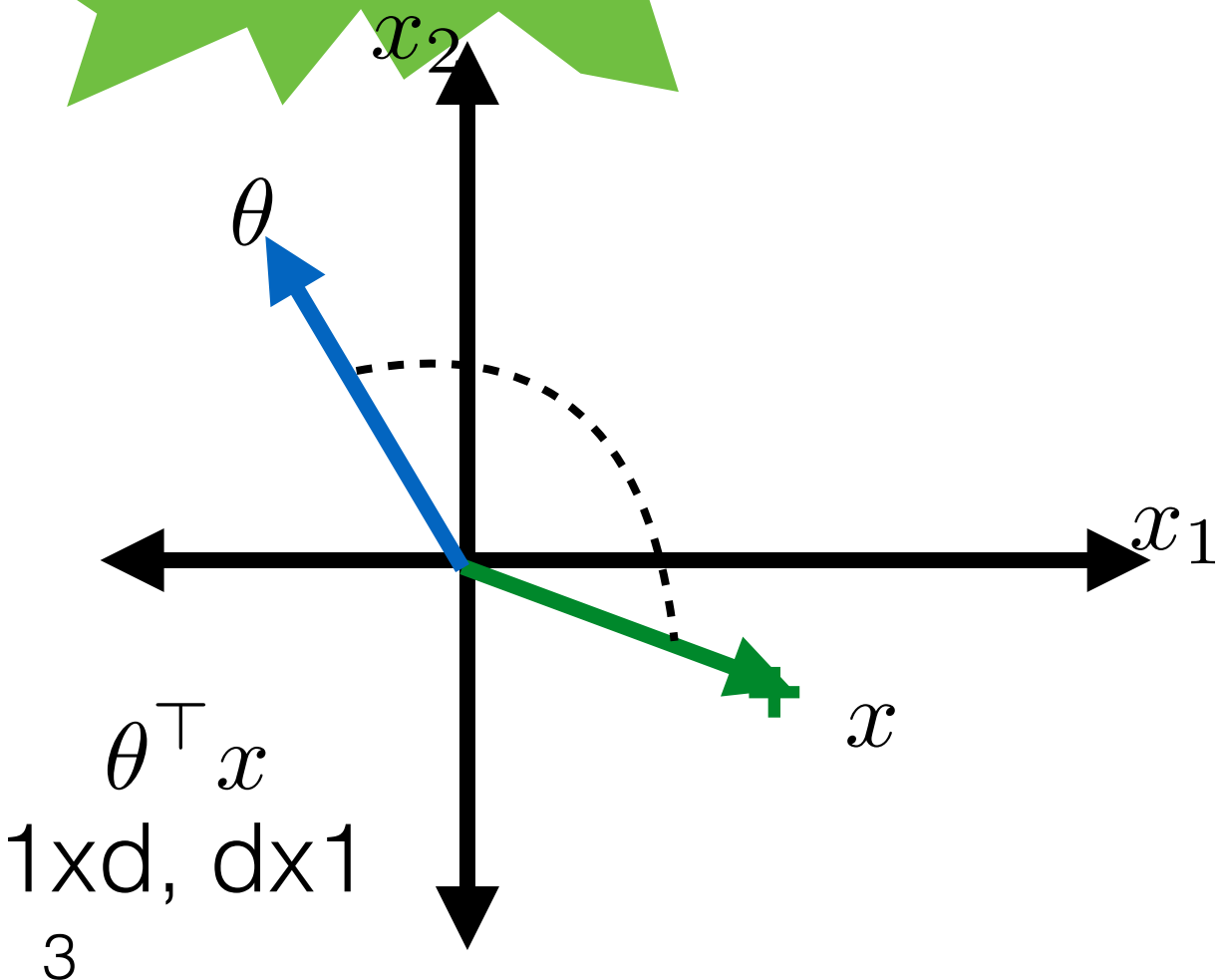Wind speed (kph)

Temperature (C)   $x_1$

3

# Linear classifiers

- Classification hypothesis:
  $$h : \mathbb{R}^d \to \{-1, +1\}$$

- Linear classifiers $\mathcal{H}$: Hypotheses that label +1 on one side of a line & -1 on the other side

$y$ = Wearing a coat?



3

# Linear classifiers

- Classification hypothesis:
  $$h : \mathbb{R}^d \to \{-1, +1\}$$

- Linear classifiers $\mathcal{H}$: Hypotheses that label +1 on one side of a line & -1 on the other side

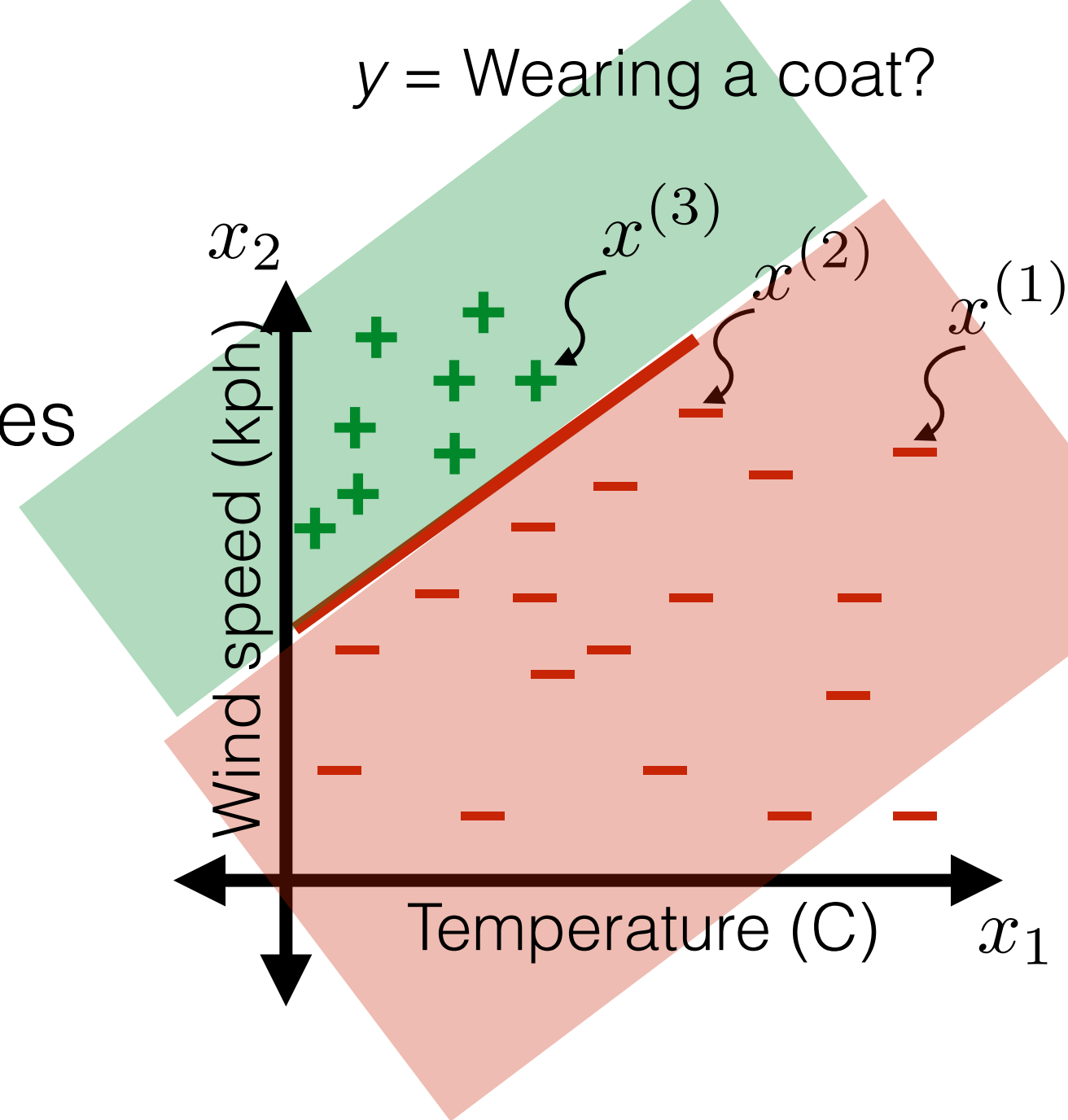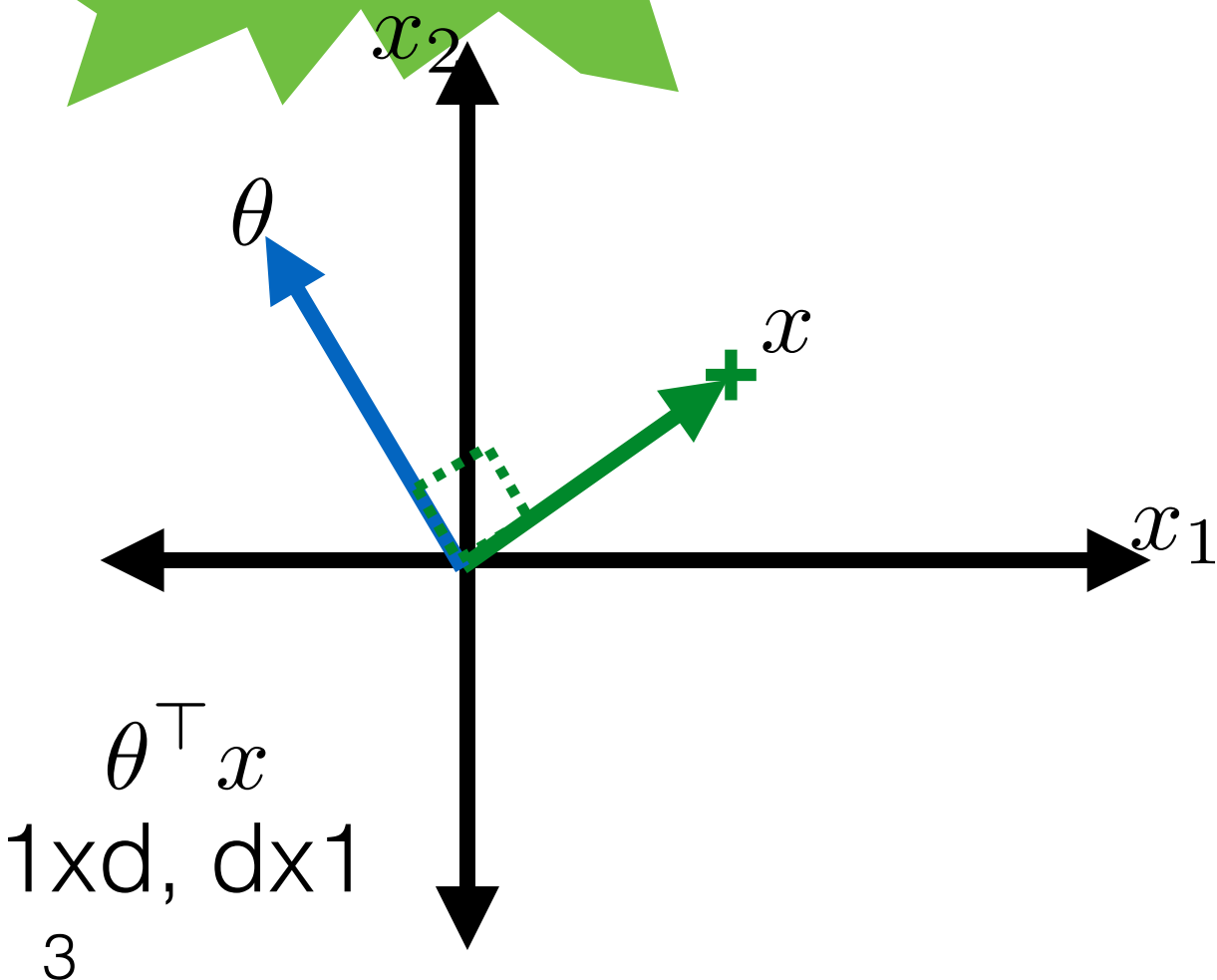**Math facts!**

$x_2$

$\theta$

$x_1$

$x$

$\theta^\top x$
1xd, dx1

3

$y$ = Wearing a coat?

$x_2$

$x^{(3)}$

$x^{(2)}$

$x^{(1)}$

Wind speed (kph)

Temperature (C)

$x_1$

# Linear classifiers

- Classification hypothesis:
  $h : \mathbb{R}^d \to \{-1, +1\}$

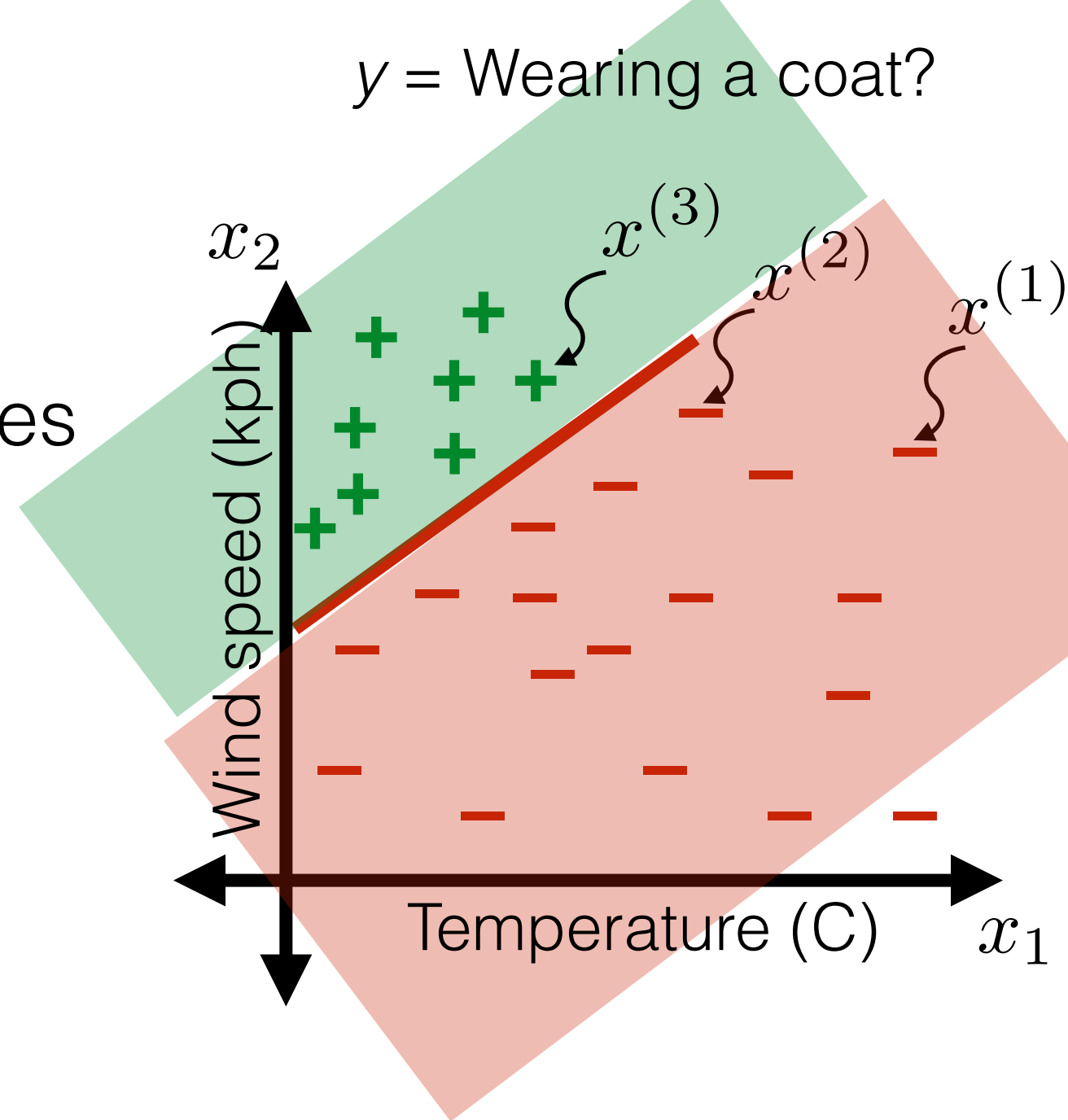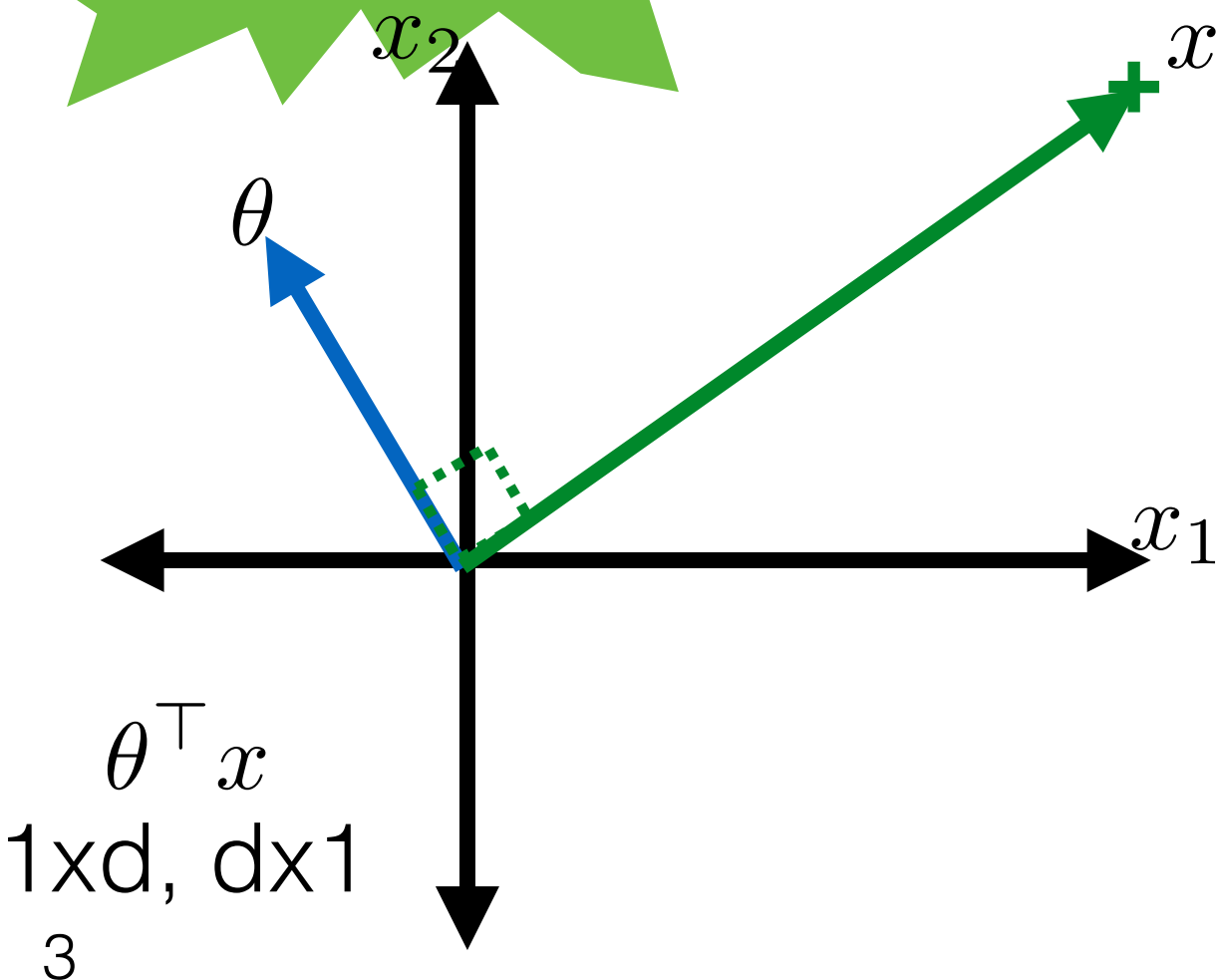- Linear classifiers $\mathcal{H}$: Hypotheses that label +1 on one side of a line & -1 on the other side

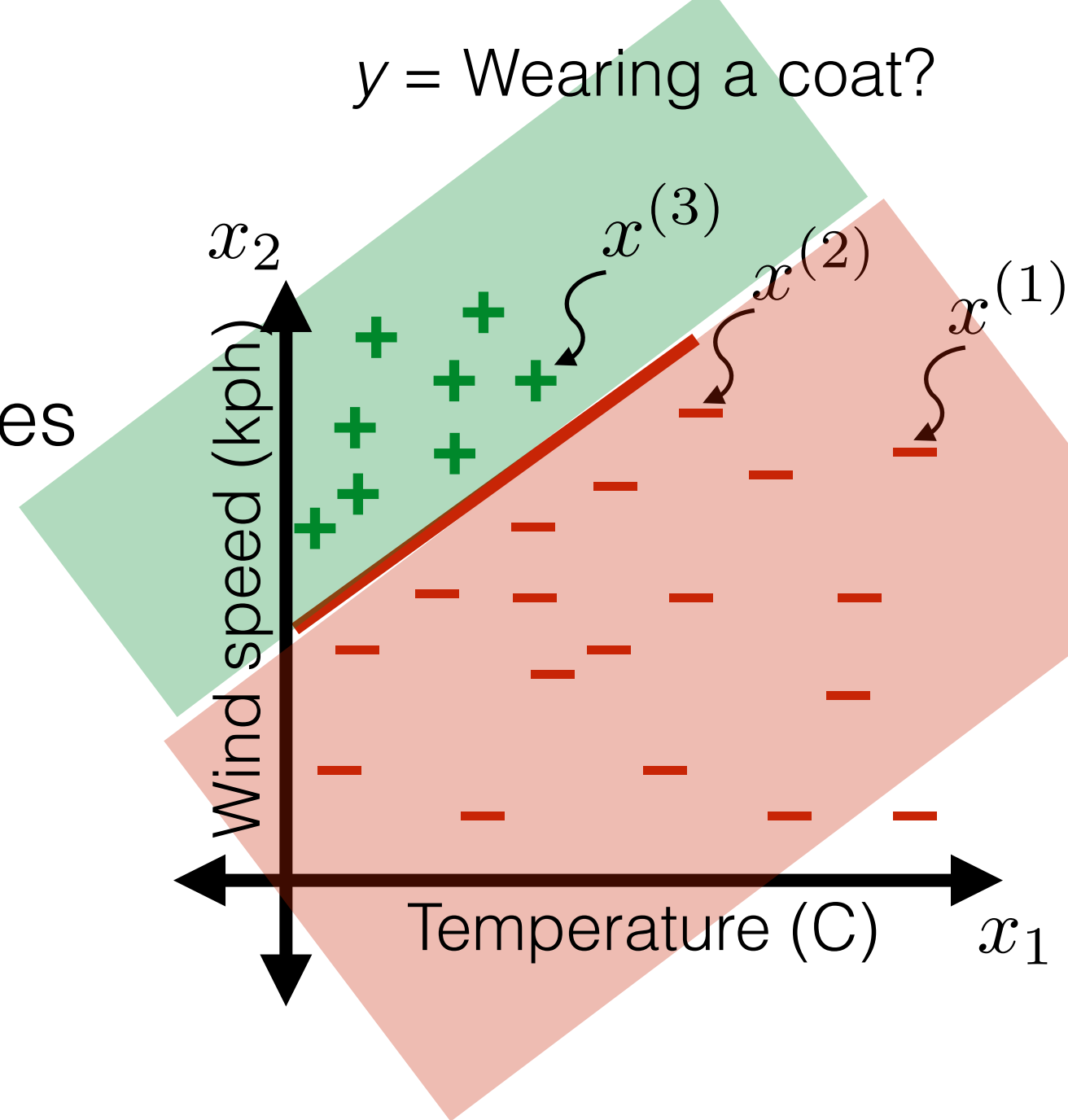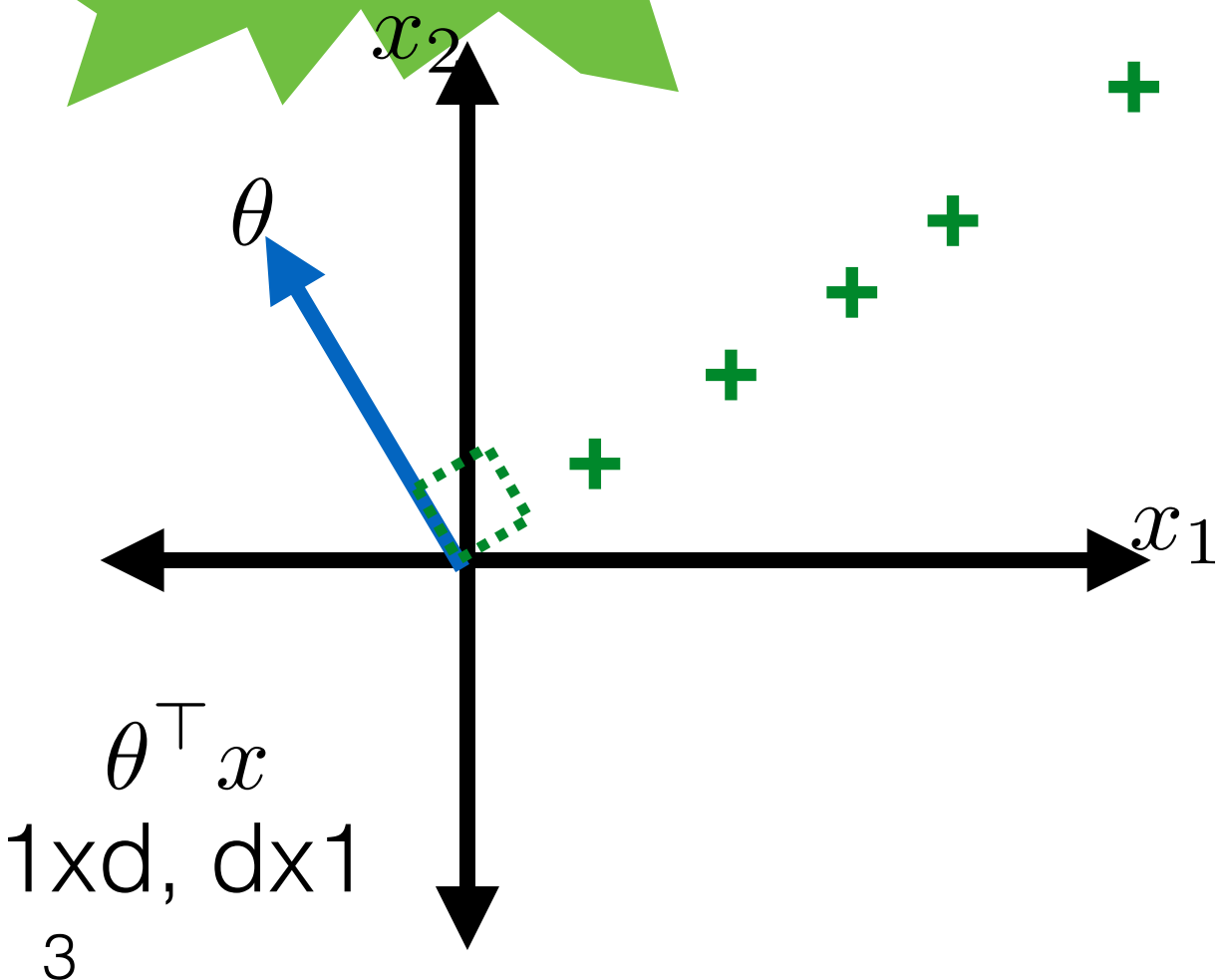**Math facts!**

$x_2$

$\theta$

$x$

$x_1$

$\theta^\top x$

1xd, dx1

3

$y$ = Wearing a coat?

$x_2$

$x^{(3)}$

$x^{(2)}$

$x^{(1)}$

Wind speed (kph)

Temperature (C)

$x_1$

# Linear classifiers

- Classification hypothesis:
  $$h : \mathbb{R}^d \to \{-1, +1\}$$

- Linear classifiers $\mathcal{H}$: Hypotheses that label +1 on one side of a line & -1 on the other side

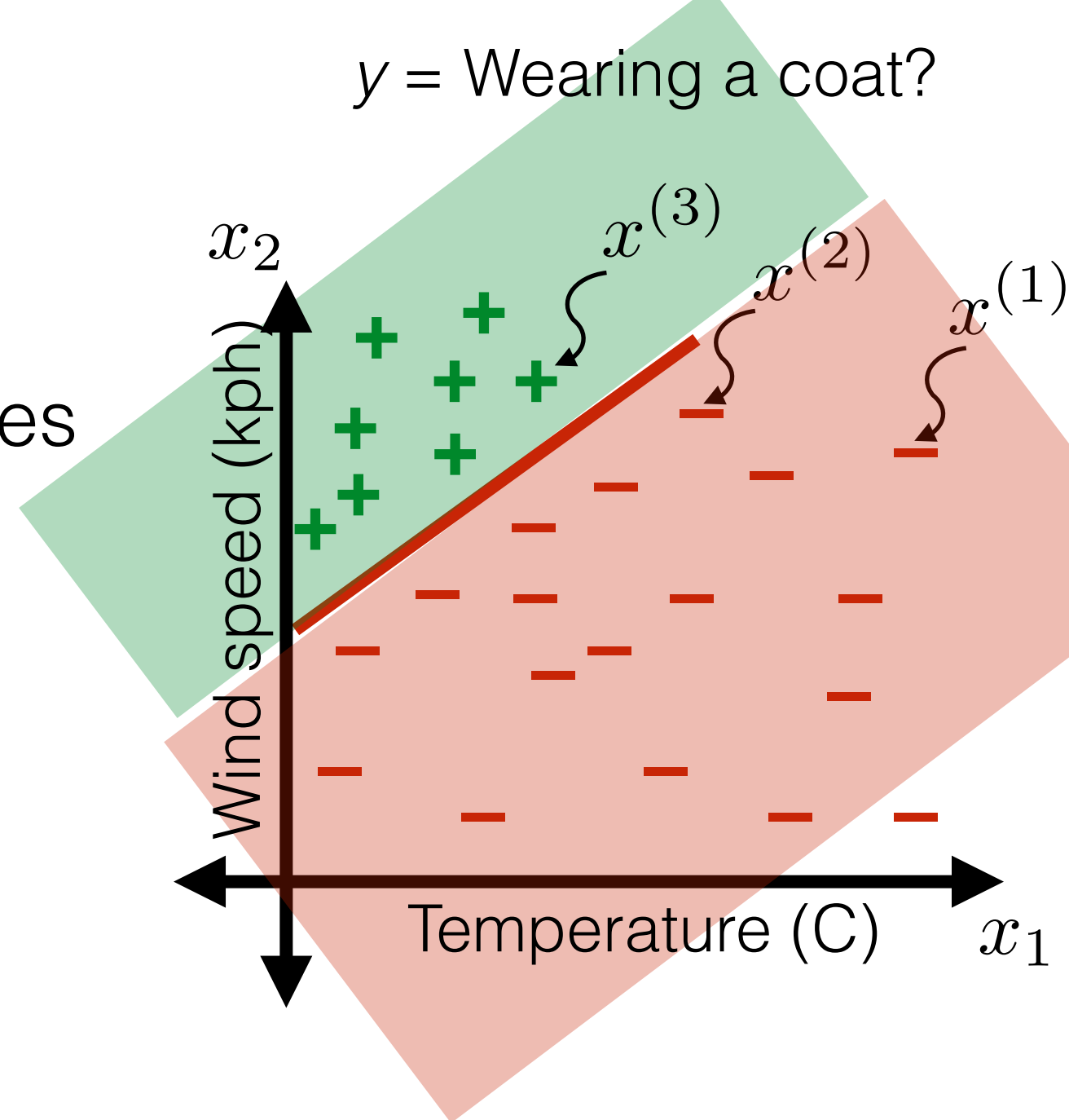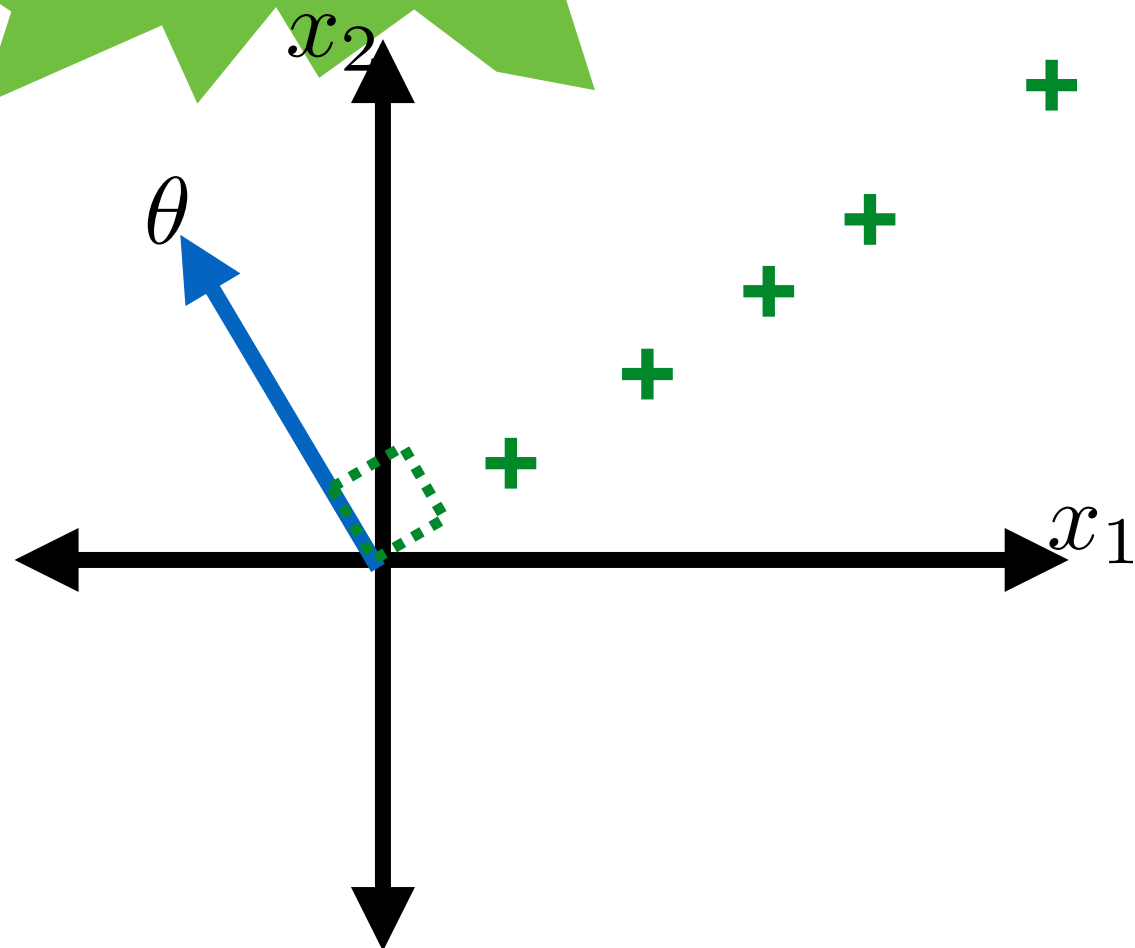**Math facts!**

$x_2$

$\theta$

$x$

$x_1$

$\theta^\top x$
1xd, dx1

3

$y$ = Wearing a coat?

$x_2$

$x^{(3)}$

$x^{(2)}$

$x^{(1)}$

Wind speed (kph)

Temperature (C)

$x_1$

# Linear classifiers

- Classification hypothesis:
  $$h : \mathbb{R}^d \to \{-1, +1\}$$

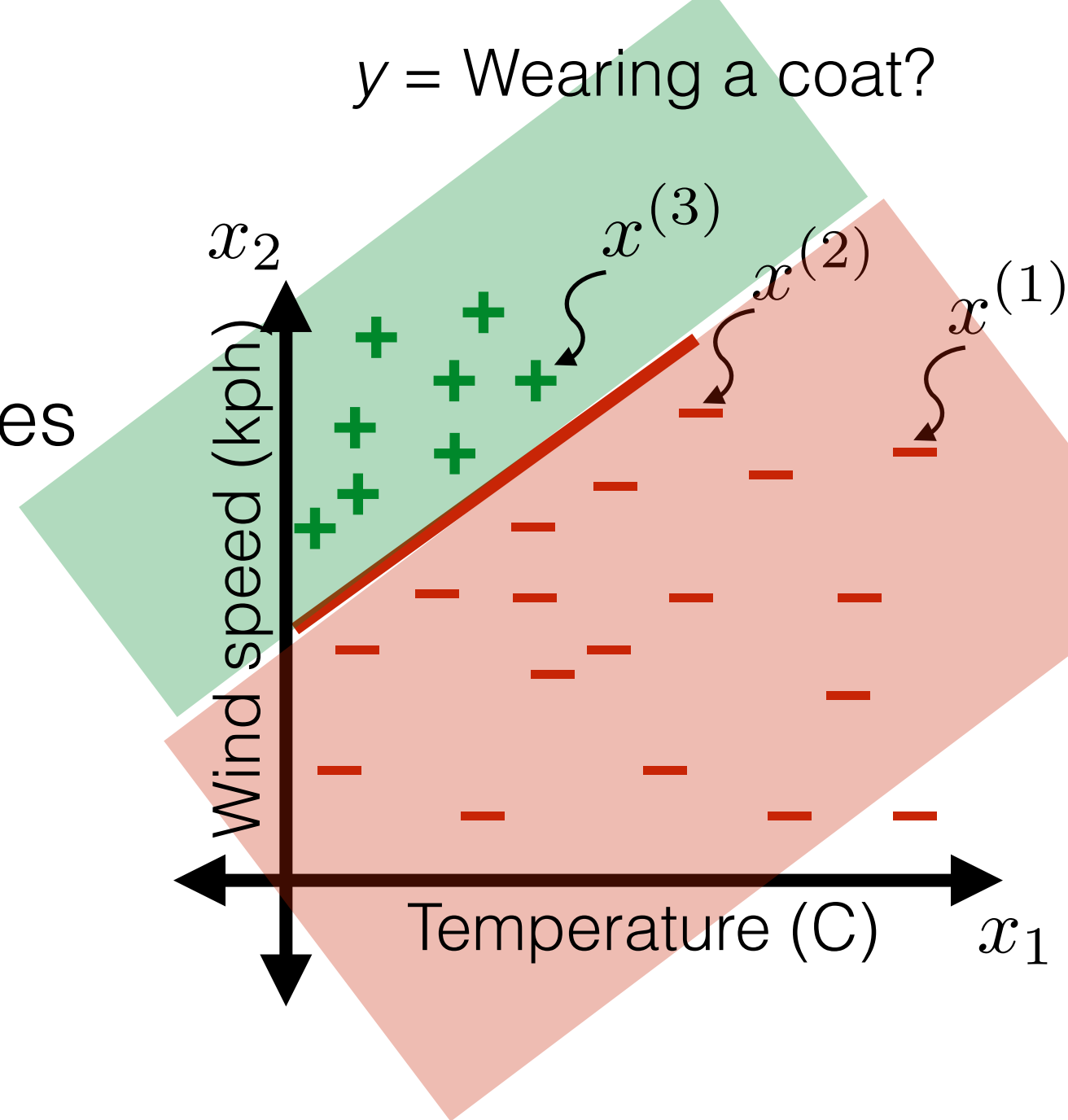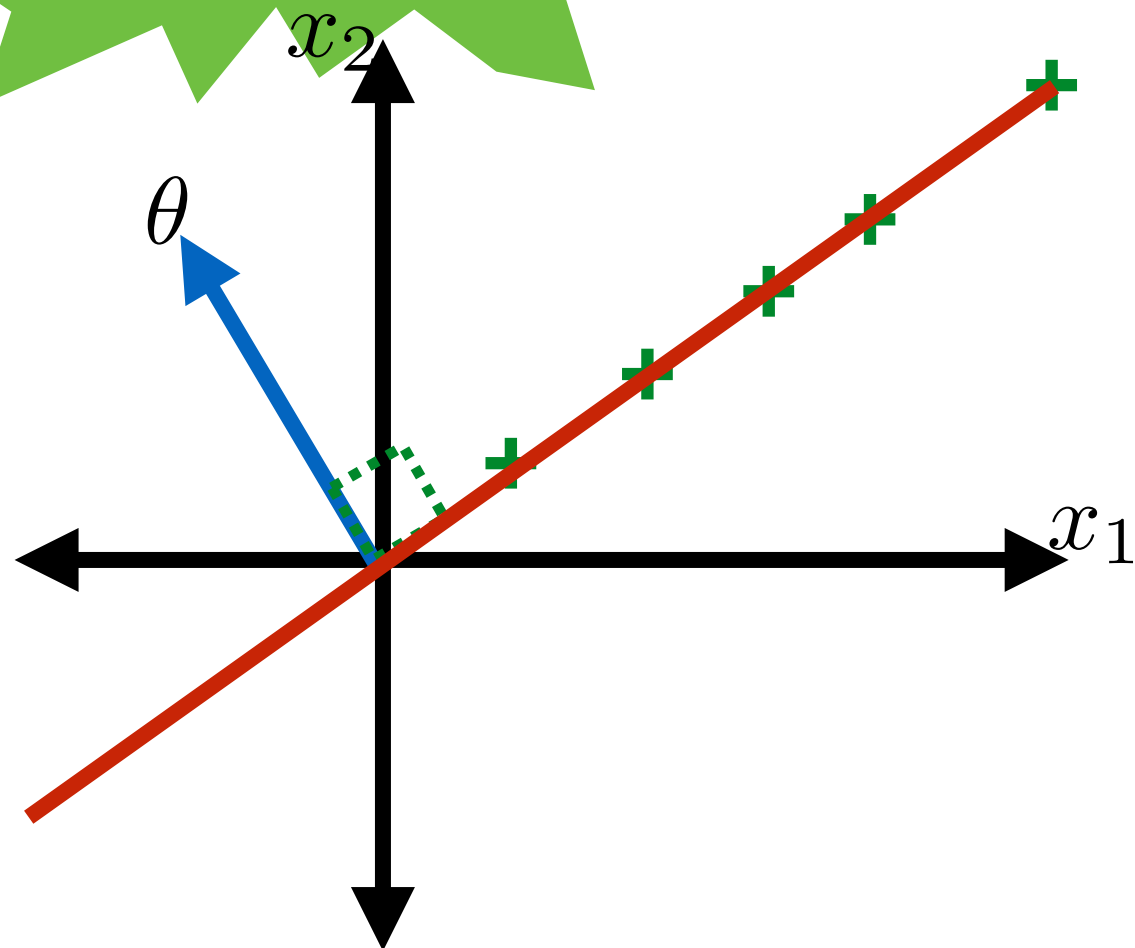- Linear classifiers $\mathcal{H}$: Hypotheses that label +1 on one side of a line & -1 on the other side

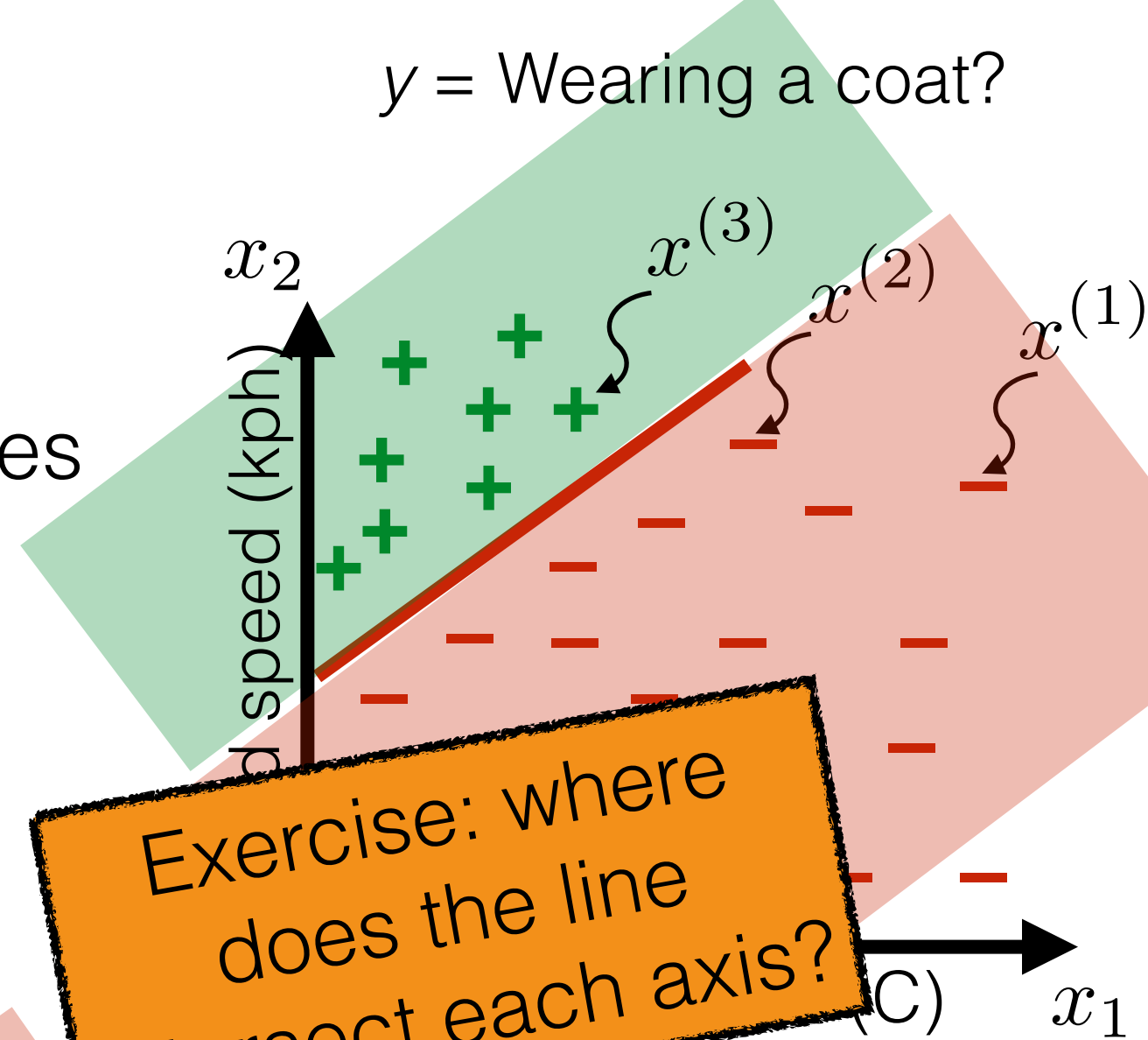**Math facts!**

$x_2$

$\theta$

$x_1$

$\theta^\top x$
1xd, dx1

3

$y$ = Wearing a coat?

$x_2$

$x^{(3)}$

$x^{(2)}$

$x^{(1)}$

Wind speed (kph)

Temperature (C)

$x_1$

# Linear classifiers

- Classification hypothesis:
  $$h : \mathbb{R}^d \to \{-1, +1\}$$

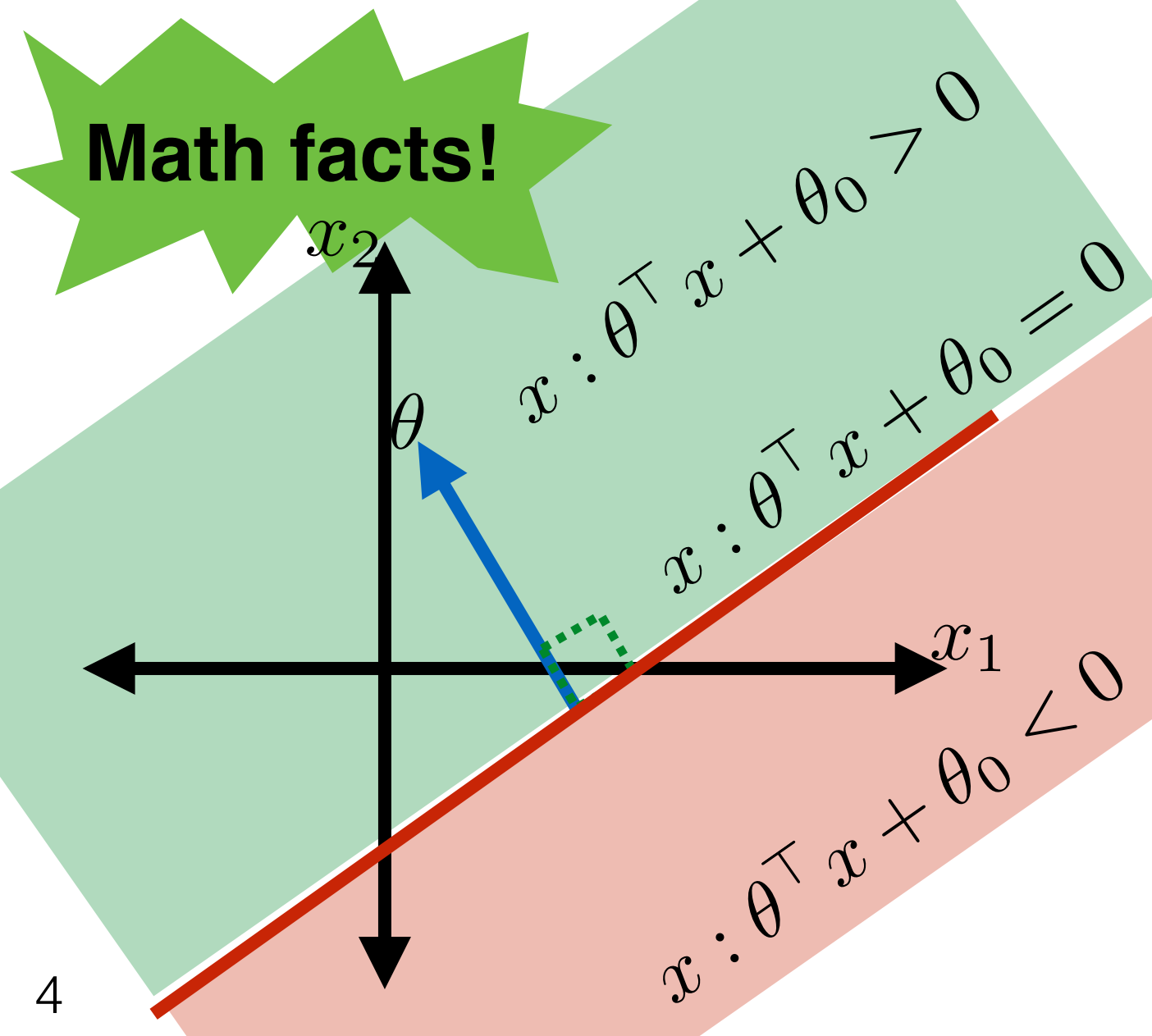- Linear classifiers $\mathcal{H}$: Hypotheses that label +1 on one side of a line & -1 on the other side
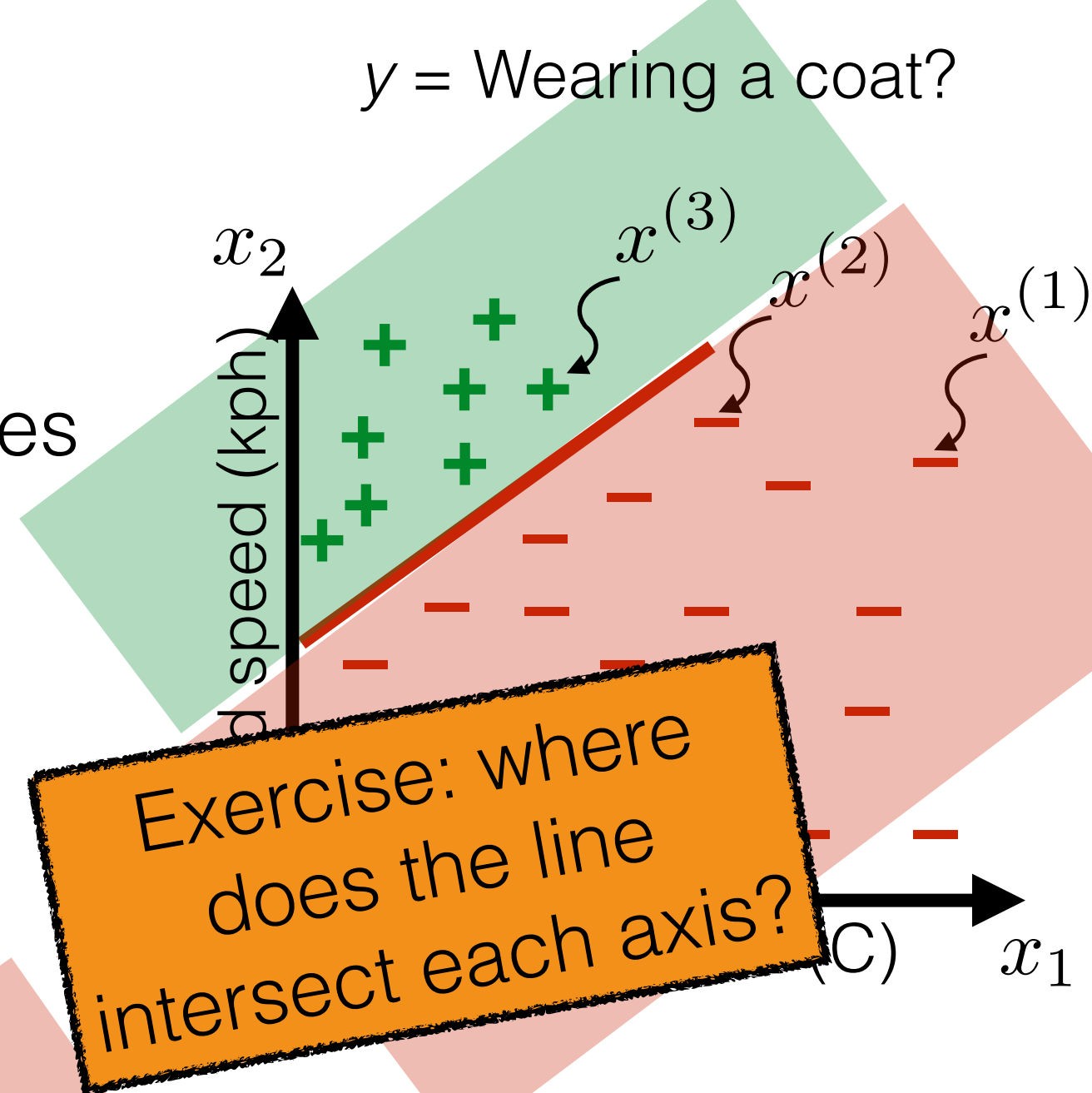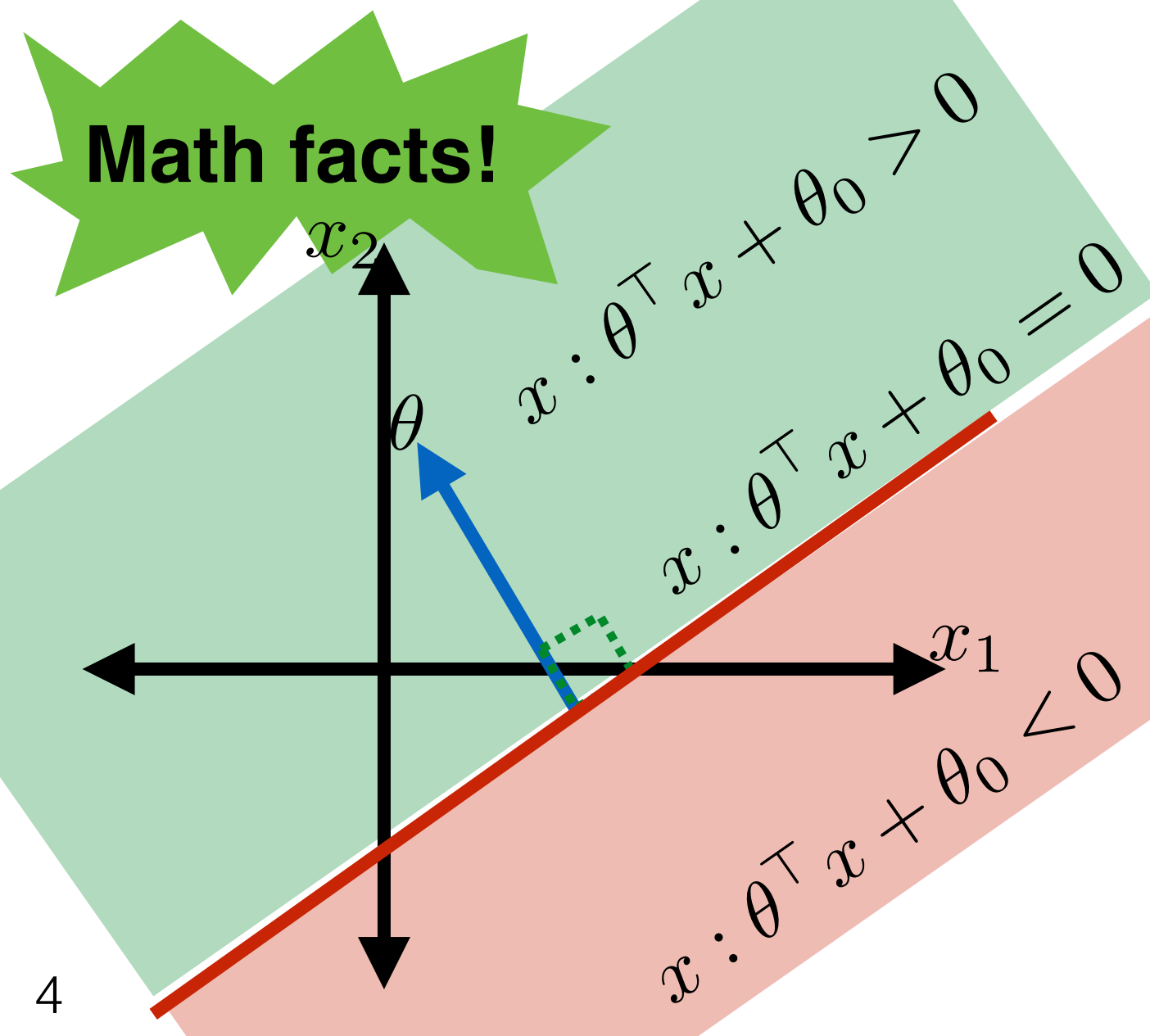
**Math facts!**

$x_2$

$\theta$

$x_1$

$y$ = Wearing a coat?

$x_2$

$x^{(3)}$

$x^{(2)}$

$x^{(1)}$

Wind speed (kph)

Temperature (C)

$x_1$

# Linear classifiers

- Classification hypothesis:
  $$h : \mathbb{R}^d \to \{-1, +1\}$$

- Linear classifiers $\mathcal{H}$: Hypotheses that label +1 on one side of a line & -1 on the other side
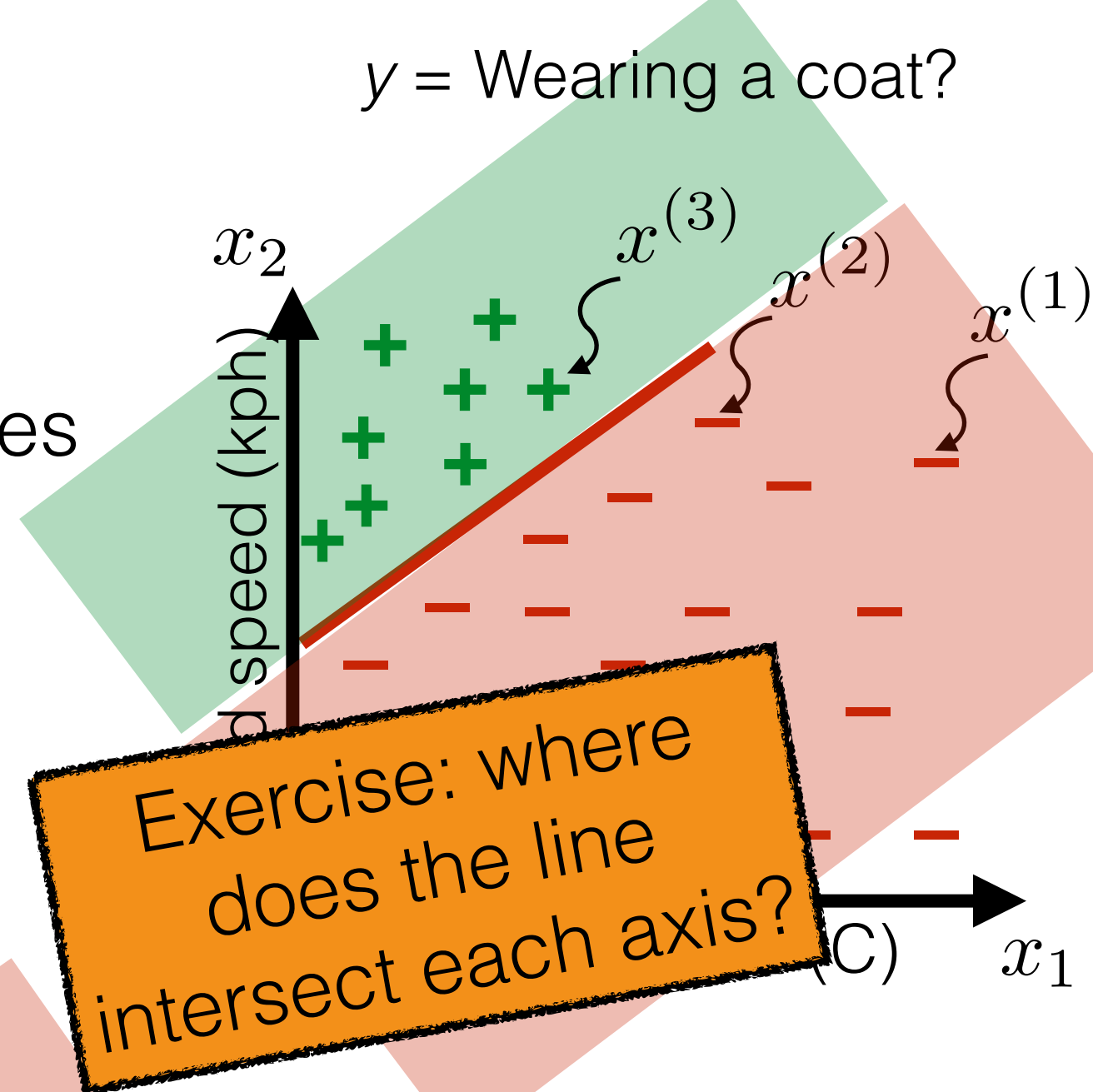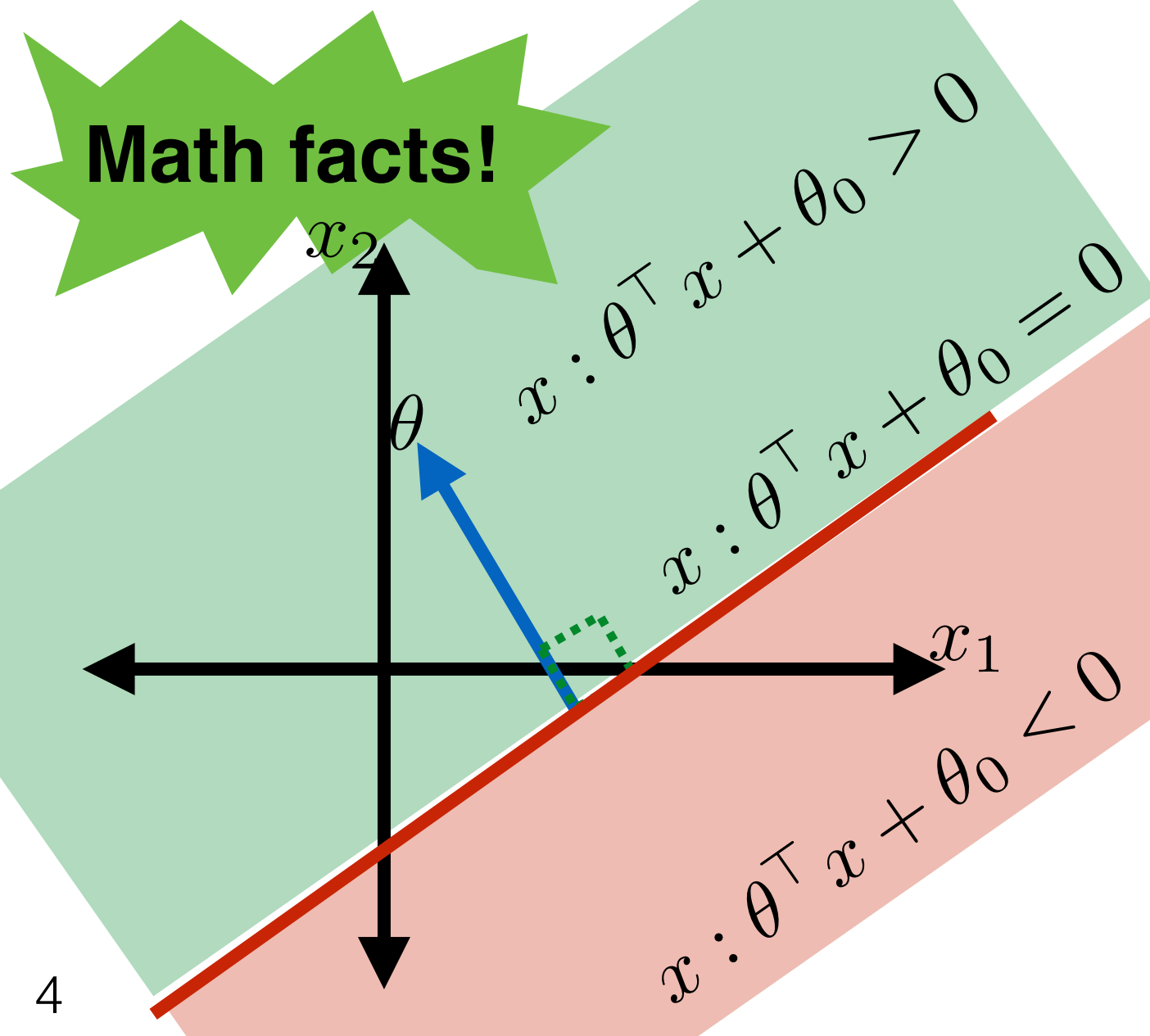
**Math facts!**

$y$ = Wearing a coat?

# Linear classifiers

- Classification hypothesis:
  $h : \mathbb{R}^d \to \{-1, +1\}$

- Linear classifiers $\mathcal{H}$: Hypotheses that label +1 on one side of a line & -1 on the other side
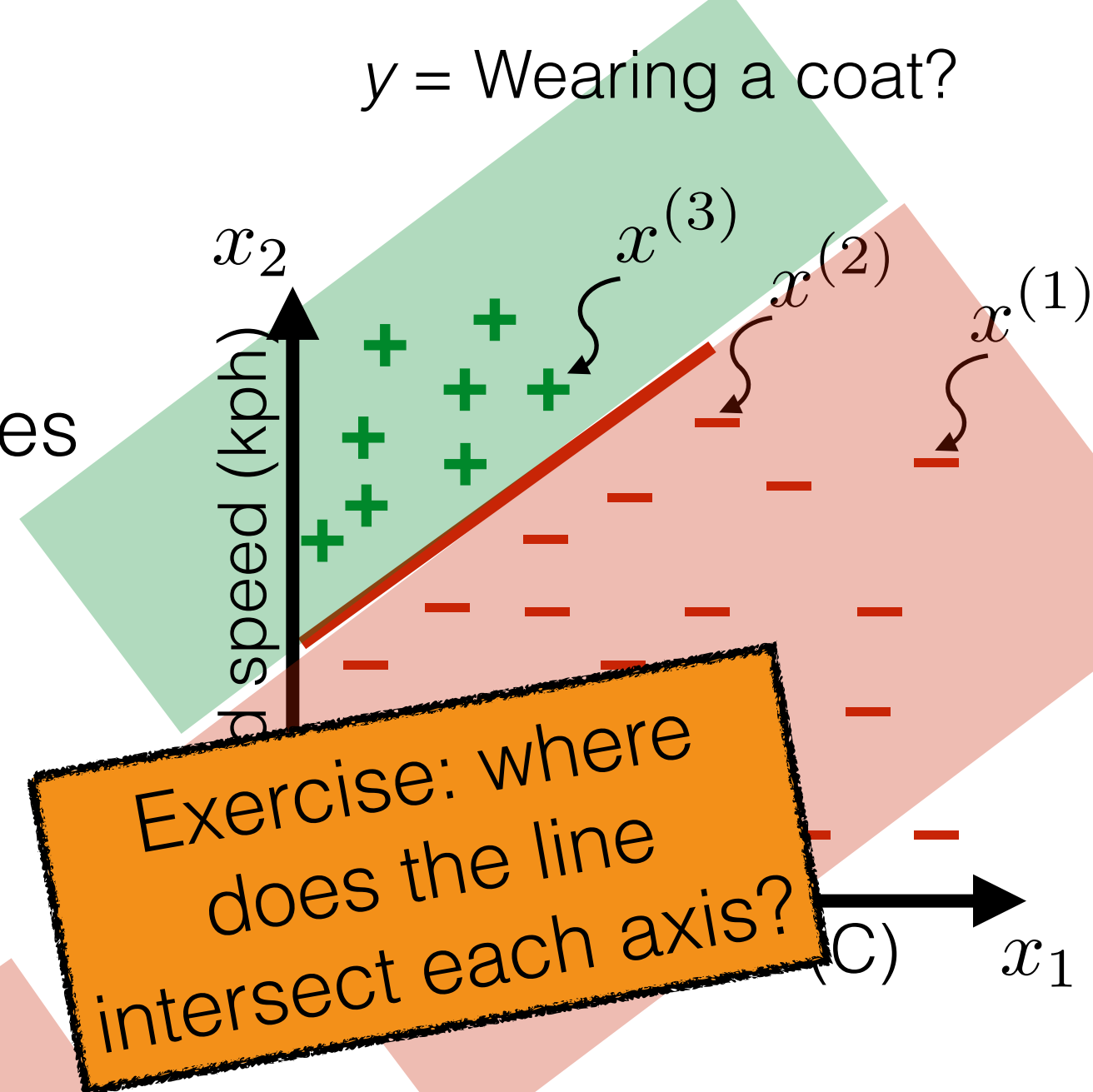
**Math facts!**

$x : \theta^\top x + \theta_0 > 0$

$x : \theta^\top x + \theta_0 = 0$

$x : \theta^\top x + \theta_0 < 0$

$x_2$

$\theta$

$x_1$

$y$ = Wearing a coat?

$x_2$

speed (kph)

$x^{(3)}$

$x^{(2)}$

$x^{(1)}$

Exercise: where does the line intersect each axis?

(C)

$x_1$

- Linear classifier:

4

# Linear classifiers

- Classification hypothesis:
  $$h : \mathbb{R}^d \to \{-1, +1\}$$

- Linear classifiers $\mathcal{H}$: Hypotheses that label +1 on one side of a line & -1 on the other side
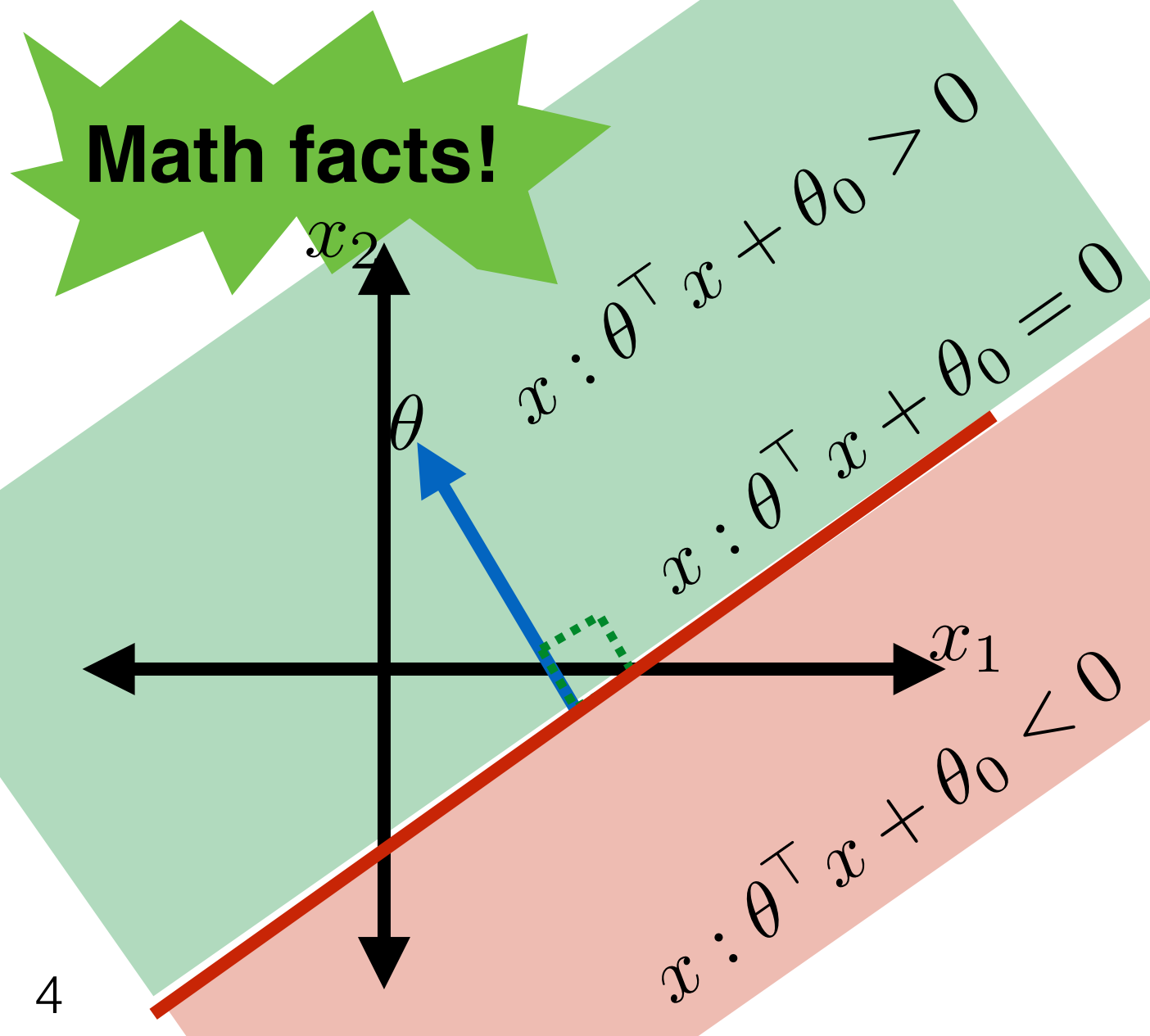
**Math facts!**

$x : \theta^\top x + \theta_0 > 0$

$x : \theta^\top x + \theta_0 = 0$

$x : \theta^\top x + \theta_0 < 0$

$x_2$

$\theta$

$x_1$

$y$ = Wearing a coat?

$x_2$

d speed (kph)

$x^{(3)}$

$x^{(2)}$

$x^{(1)}$

Exercise: where does the line intersect each axis?

(C)

$x_1$

- Linear classifier:
  $$h(x; \theta, \theta_0) = \text{sign}(\theta^\top x + \theta_0)$$
  $$= \begin{cases} +1 \text{ if } \theta^\top x + \theta_0 > 0 \\ -1 \text{ if } \theta^\top x + \theta_0 < 0 \end{cases}$$
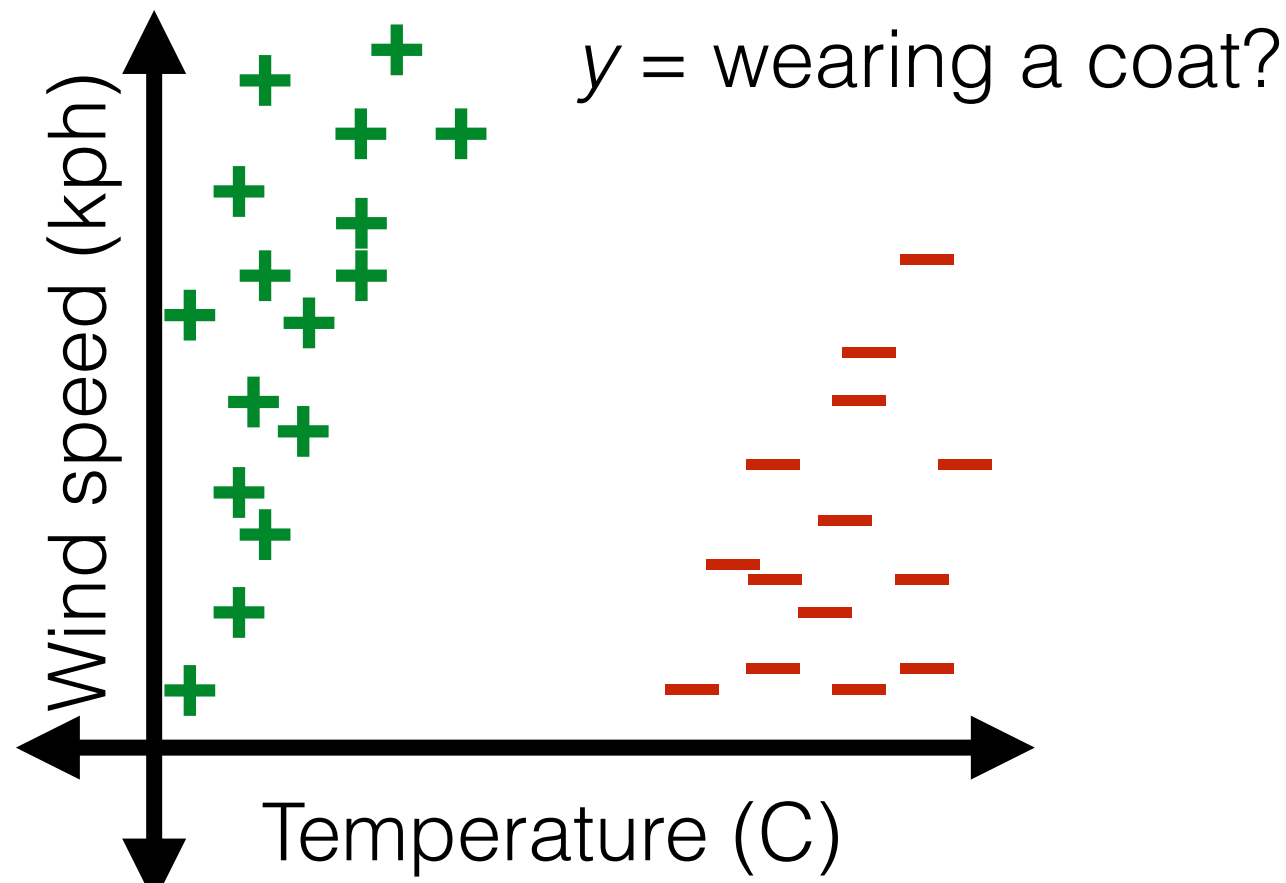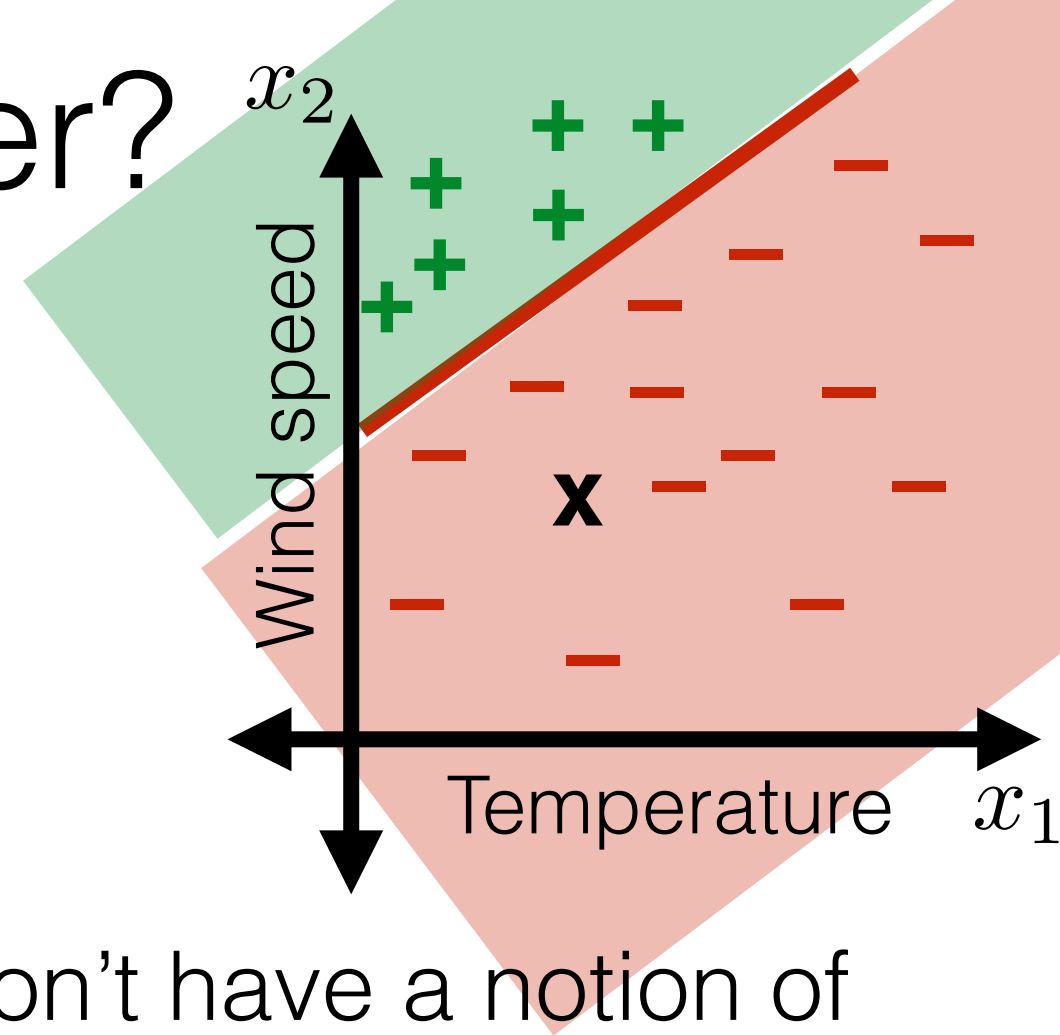
# Linear classifiers

- Classification hypothesis:
  $$h : \mathbb{R}^d \rightarrow \{-1, +1\}$$

- Linear classifiers $\mathcal{H}$: Hypotheses that label +1 on one side of a line & -1 on the other side

**Math facts!**

$x_2$

$\theta$

$x : \theta^\top x + \theta_0 > 0$

$x : \theta^\top x + \theta_0 = 0$

$x : \theta^\top x + \theta_0 < 0$

$x_1$

$y$ = Wearing a coat?

$x_2$

speed (kph)

$x^{(3)}$

$x^{(2)}$

$x^{(1)}$

$+$

$-$

Exercise: where does the line intersect each axis?

(C)

$x_1$

- Linear classifier:
$$h(x; \theta, \theta_0) = \text{sign}(\theta^\top x + \theta_0)$$

$$= \begin{cases} +1 \text{ if } \theta^\top x + \theta_0 > 0 \\ -1 \text{ if } \theta^\top x + \theta_0 \leq 0 \end{cases}$$

# Linear classifiers

- Classification hypothesis:
  $$h : \mathbb{R}^d \to \{-1, +1\}$$

- Linear classifiers $\mathcal{H}$: Hypotheses that label +1 on one side of a line & -1 on the other side

**Math facts!**

$x_2$

$\theta$

$x : \theta^\top x + \theta_0 > 0$

$x : \theta^\top x + \theta_0 = 0$

$x_1$

$x : \theta^\top x + \theta_0 < 0$

$y$ = Wearing a coat?

$x_2$

speed (kph)

$x^{(3)}$

$x^{(2)}$

$x^{(1)}$

Exercise: where does the line intersect each axis?

(C)

$x_1$

- Linear classifier:
  $$h(x; \theta, \theta_0) = \text{sign}(\theta^\top x + \theta_0)$$
  $$= \begin{cases} +1 \text{ if } \theta^\top x + \theta_0 > 0 \\ -1 \text{ if } \theta^\top x + \theta_0 \leq 0 \end{cases}$$

- Note: $\theta$ tells us direction

4

# How good is a classifier?

- Should predict well on future data
  - Example: 0-1 loss

$$L(g, a) = \begin{cases} 0 \text{ if } g = a \\ 1 \text{ else} \end{cases}$$

g: guess,
a: actual

- Example: asymmetric loss

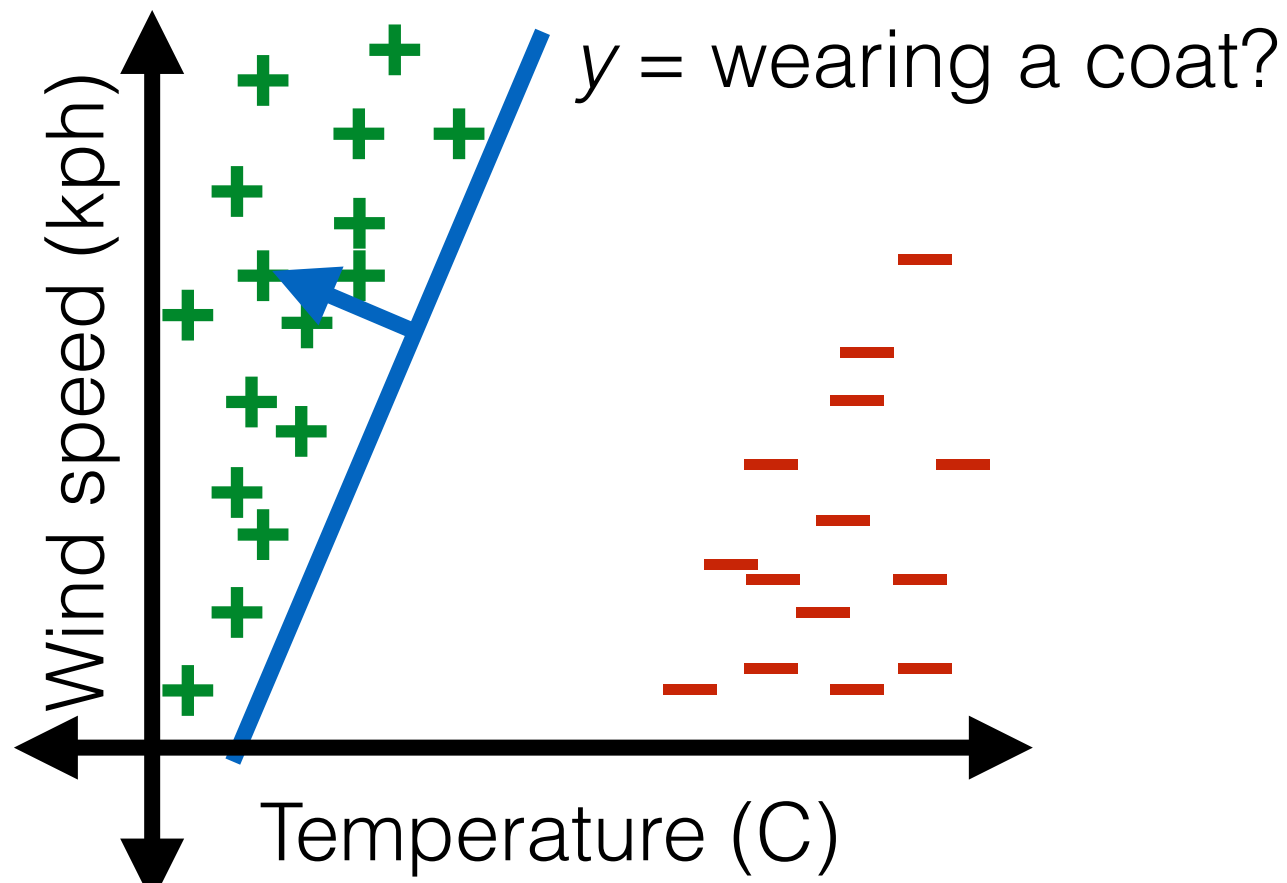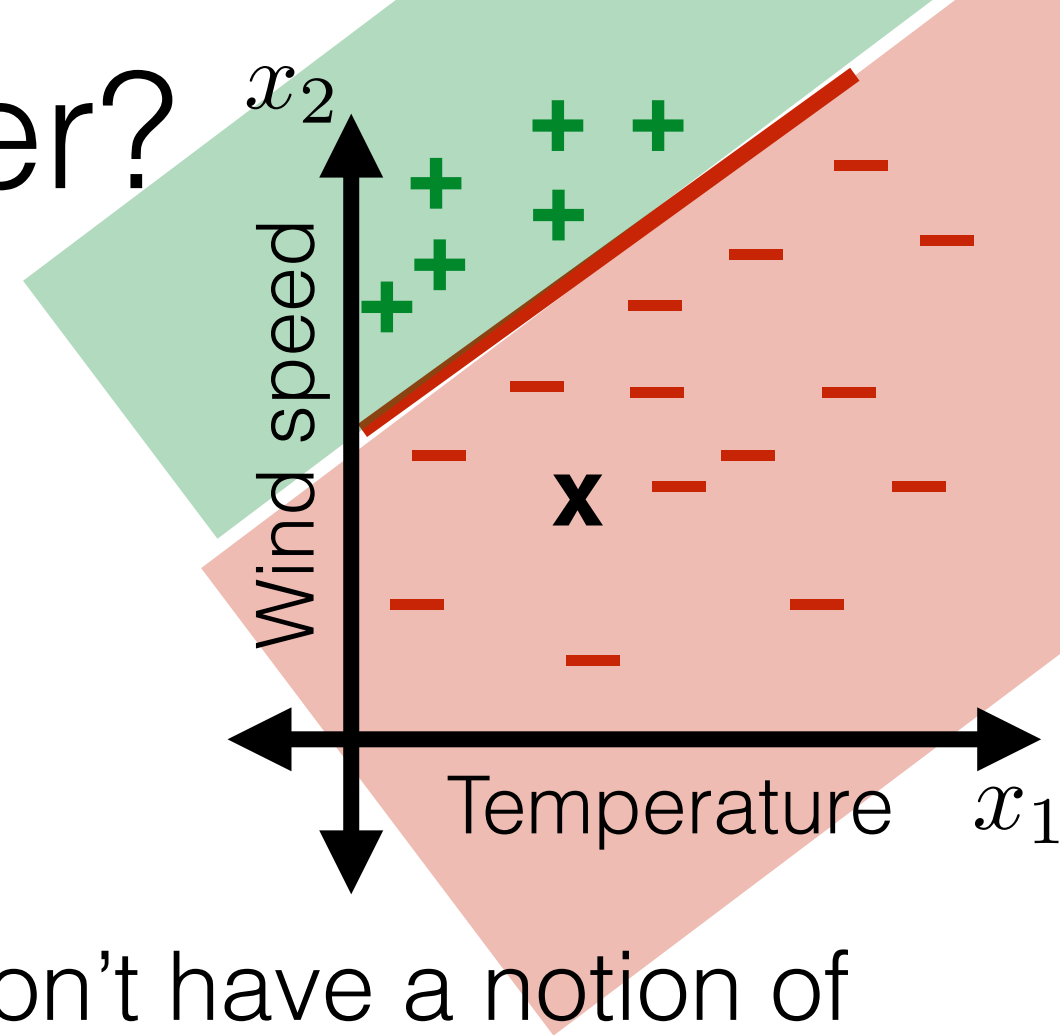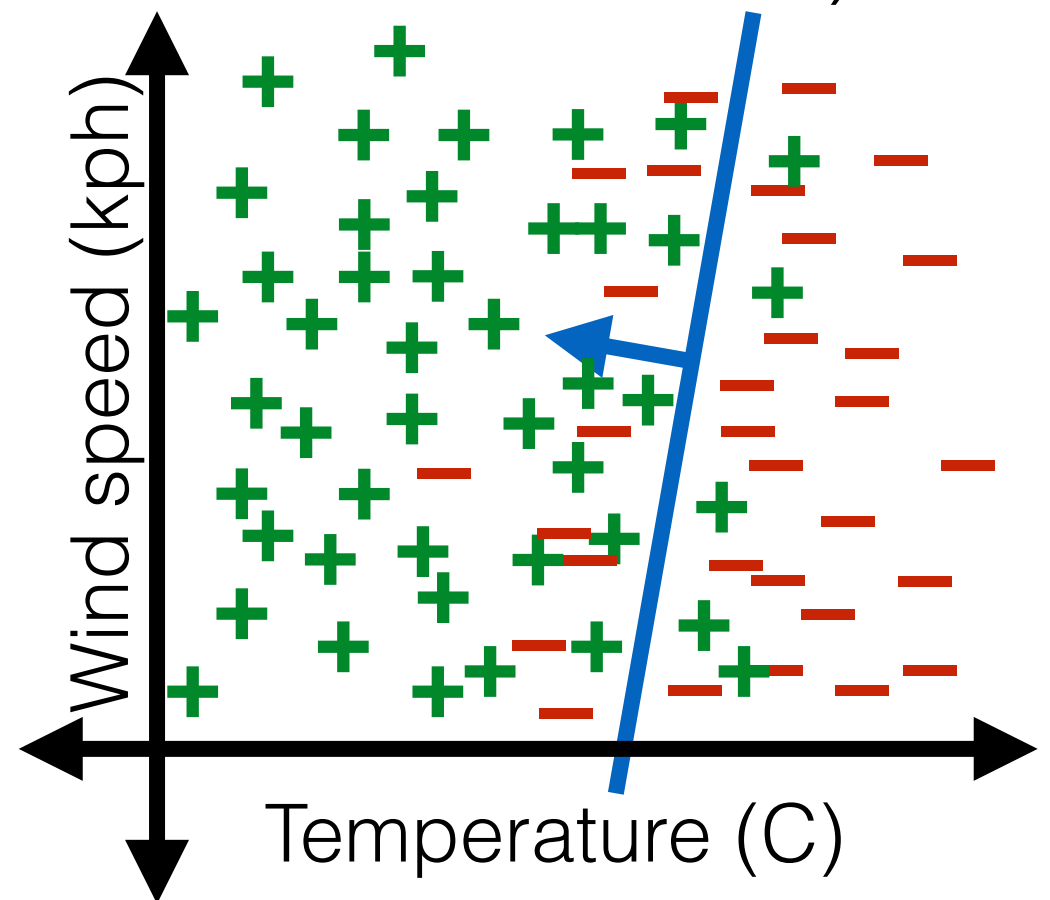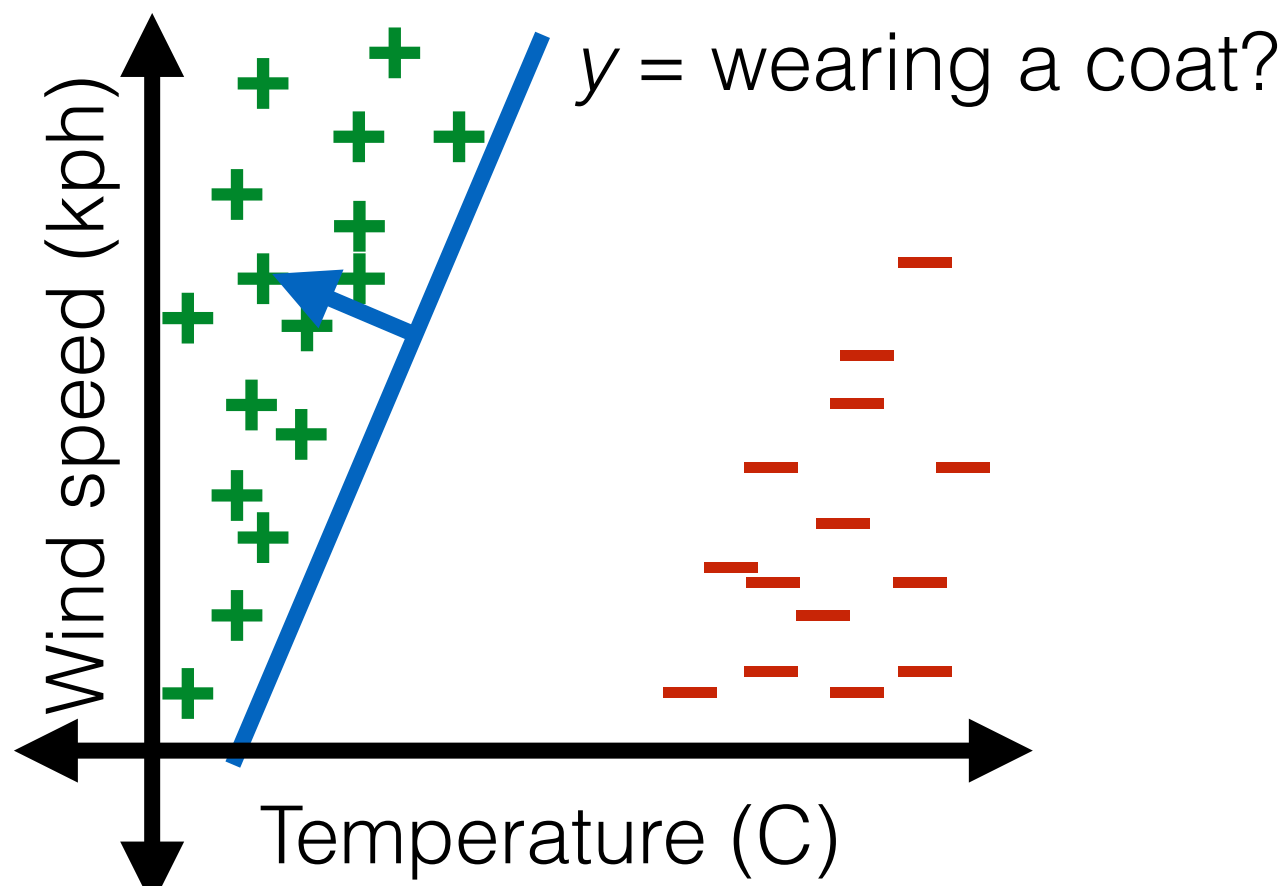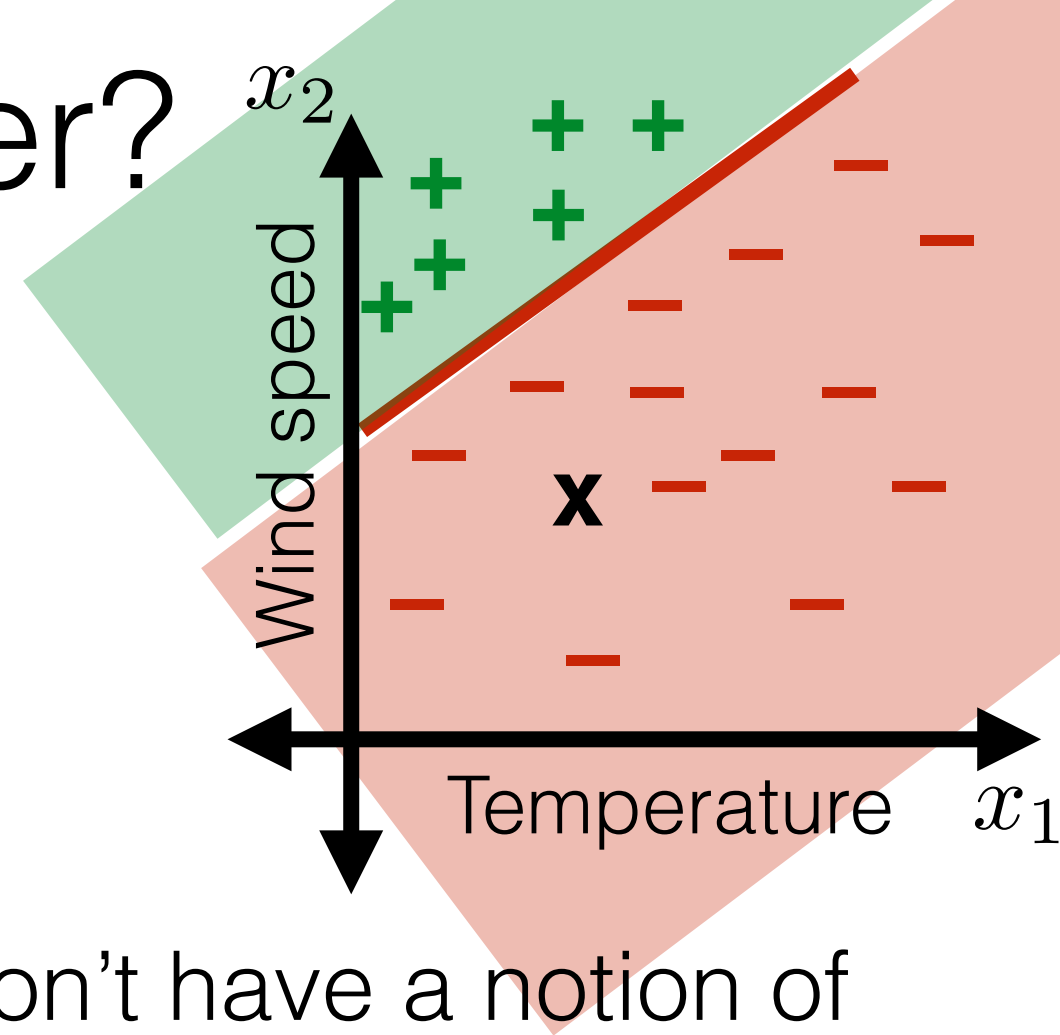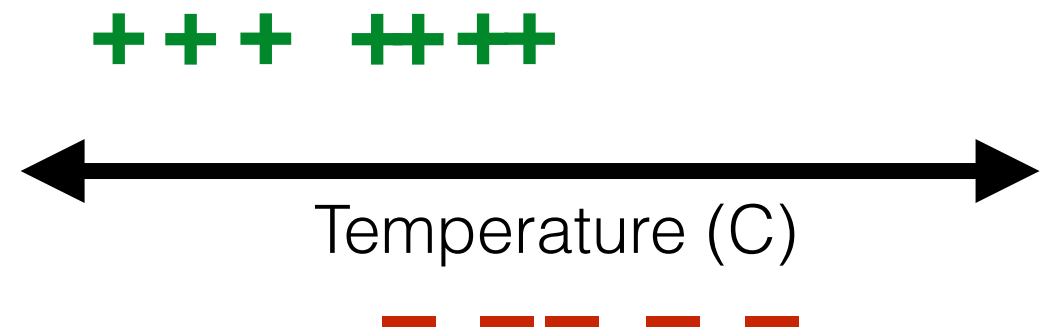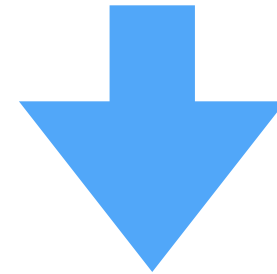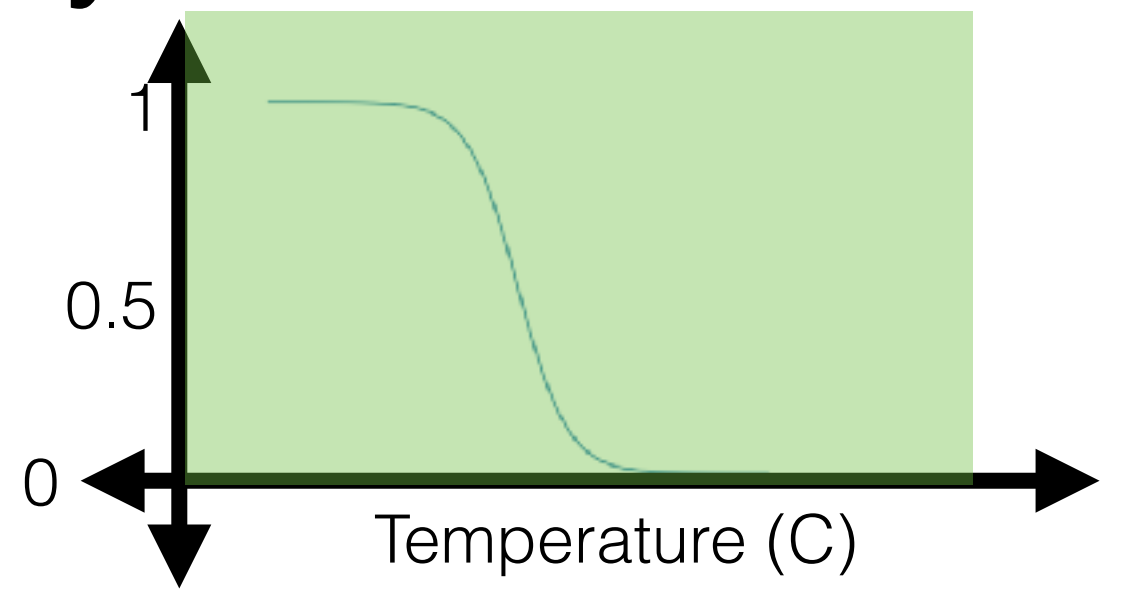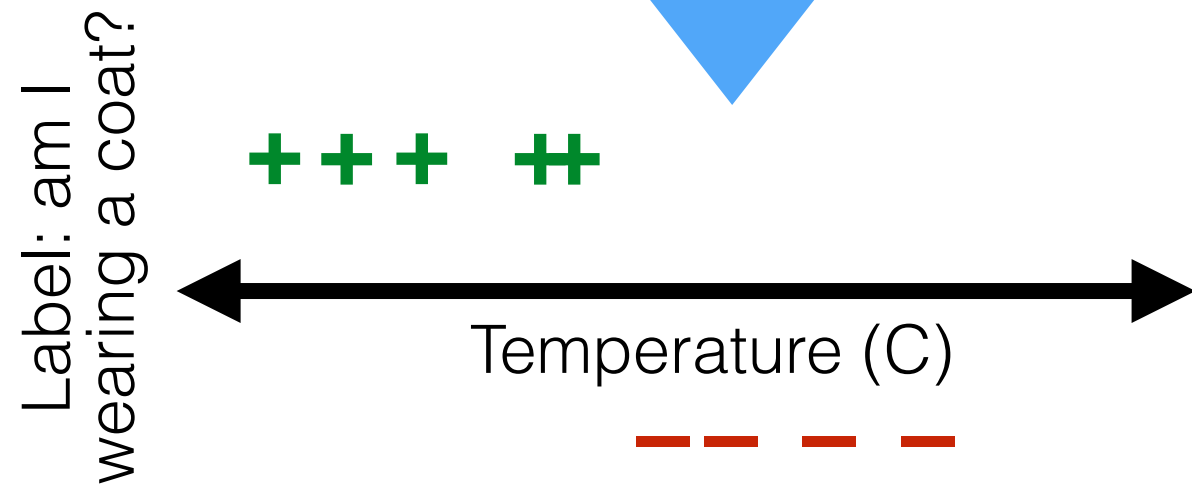- But: 0-1 loss & linear classifiers don't have a notion of uncertainty (how well do we know what we know?)

$y$ = wearing a coat?

# How good is a classifier?

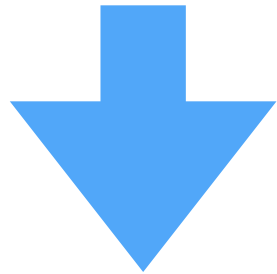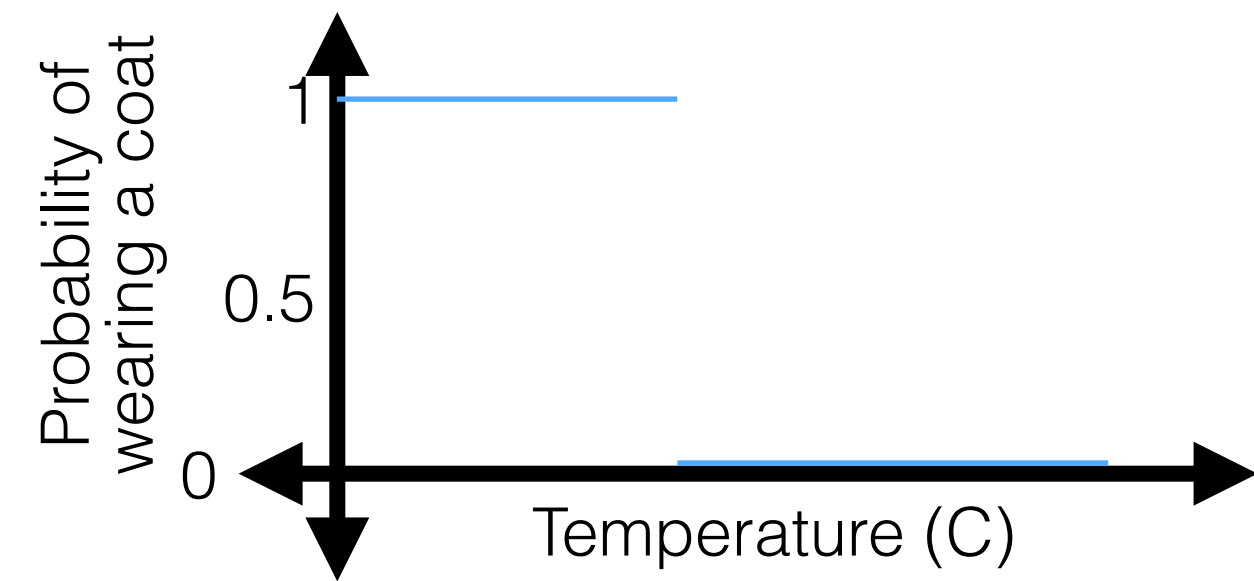- Should predict well on future data
  - Example: 0-1 loss

$$L(g,a) = \begin{cases} 0 \text{ if } g = a \\ 1 \text{ else} \end{cases}$$

g: guess,
a: actual

- Example: asymmetric loss

- But: 0-1 loss & linear classifiers don't have a notion of uncertainty (how well do we know what we know?)
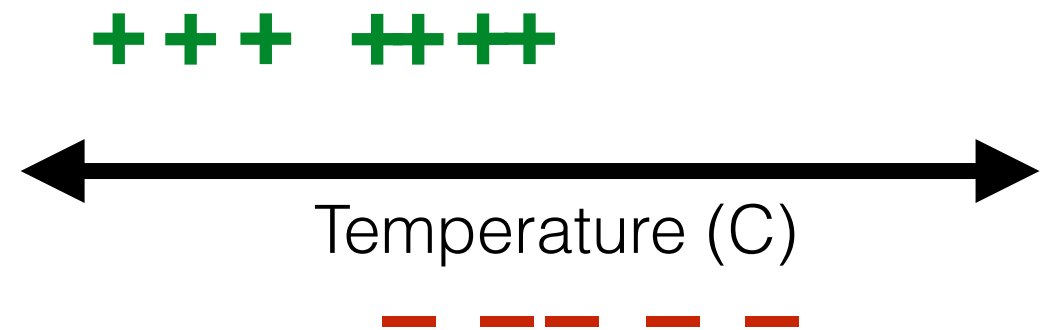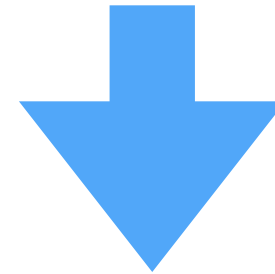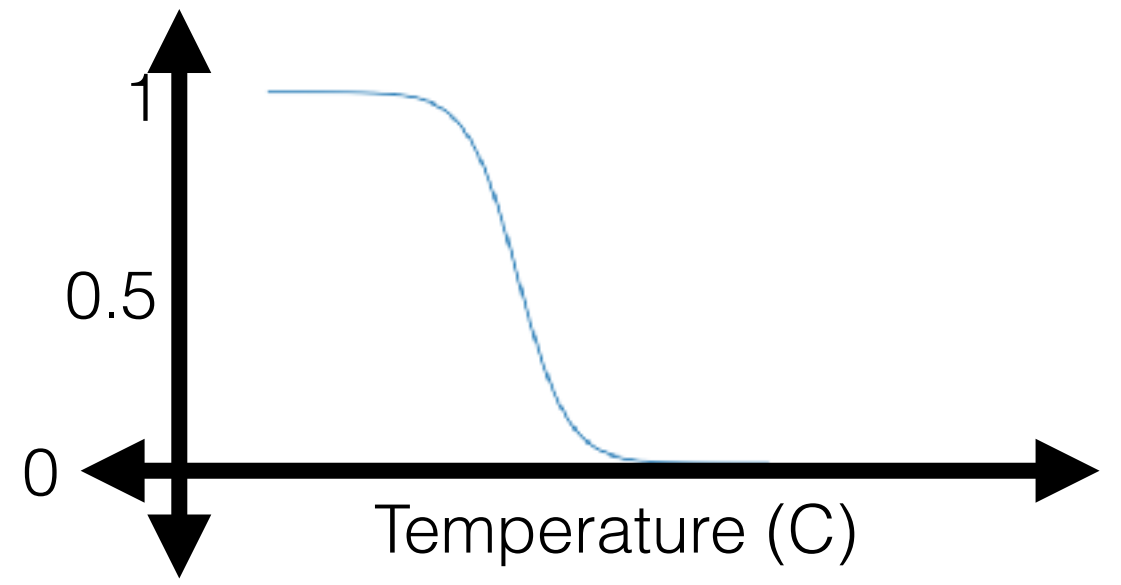
$x_2$

Wind speed

**X**

Temperature  $x_1$

$y$ = wearing a coat?

Wind speed (kph)

Temperature (C)

# How good is a classifier?

- Should predict well on future data
  - Example: 0-1 loss

$$L(g, a) = \begin{cases} 0 \text{ if } g = a \\ 1 \text{ else} \end{cases}$$

g: guess,
a: actual

  - Example: asymmetric loss

- But: 0-1 loss & linear classifiers don't have a notion of uncertainty (how well do we know what we know?)
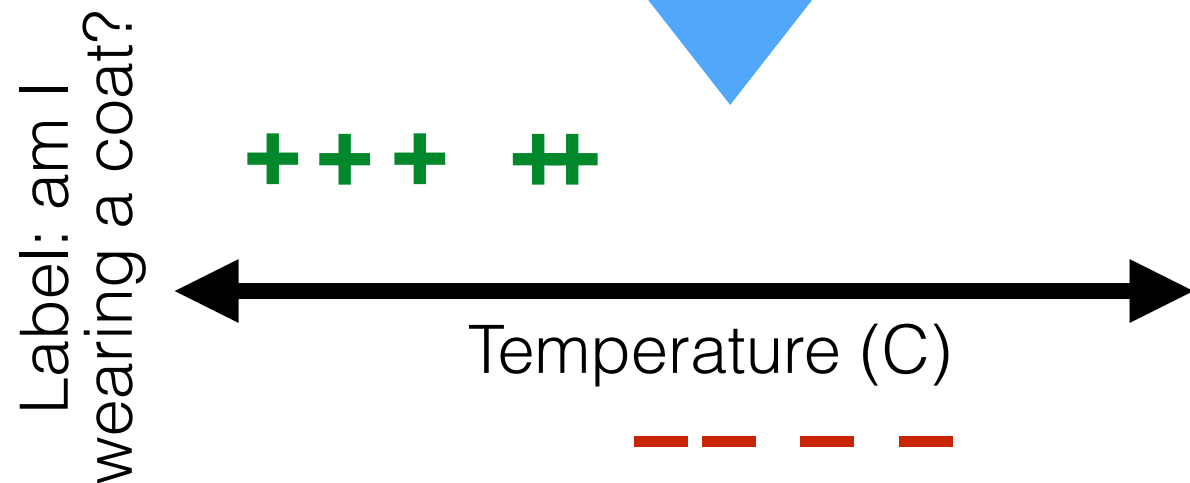
*y* = wearing a coat?

# How good is a classifier?

- Should predict well on future data
  - Example: 0-1 loss

$$L(g,a) = \begin{cases} 0 \text{ if } g = a \\ 1 \text{ else} \end{cases}$$

g: guess,
a: actual

- Example: asymmetric loss

- But: 0-1 loss & linear classifiers don't have a notion of uncertainty (how well do we know what we know?)

$y$ = wearing a coat?

# How good is a classifier?

- Should predict well on future data
  - Example: 0-1 loss

$$L(g,a) = \begin{cases} 0 \text{ if } g = a \\ 1 \text{ else} \end{cases}$$

g: guess,
a: actual

- Example: asymmetric loss

- But: 0-1 loss & linear classifiers don't have a notion of uncertainty (how well do we know what we know?)

$x_2$

Wind speed

Temperature  $x_1$

*y* = wearing a coat?

Wind speed (kph)

Temperature (C)

# How good is a classifier?

- Should predict well on future data
  - Example: 0-1 loss

$$L(g,a) = \begin{cases} 0 \text{ if } g = a \\ 1 \text{ else} \end{cases}$$

g: guess,
a: actual

  - Example: asymmetric loss

- But: 0-1 loss & linear classifiers don't have a notion of uncertainty (how well do we know what we know?)

*y* = wearing a coat?

# Capturing uncertainty



- How to make this shape?

# Capturing uncertainty



- How to make this shape?
  - Sigmoid/logistic function
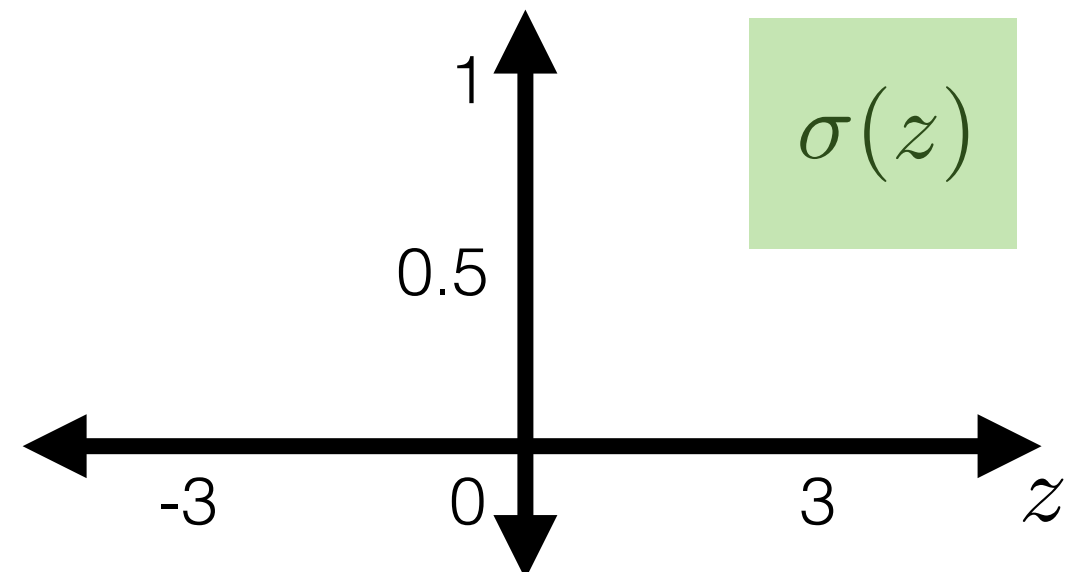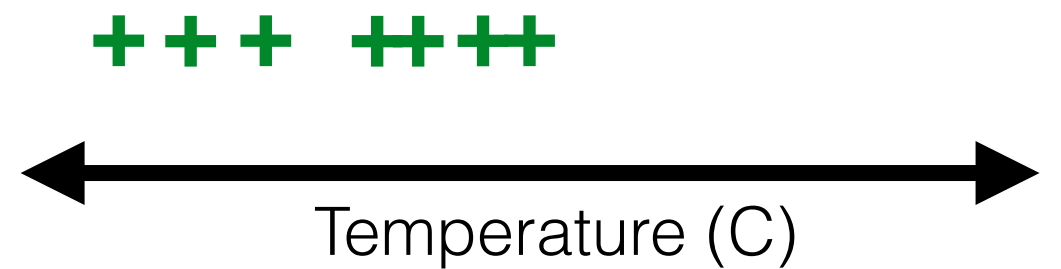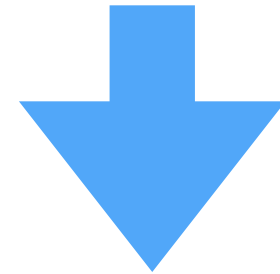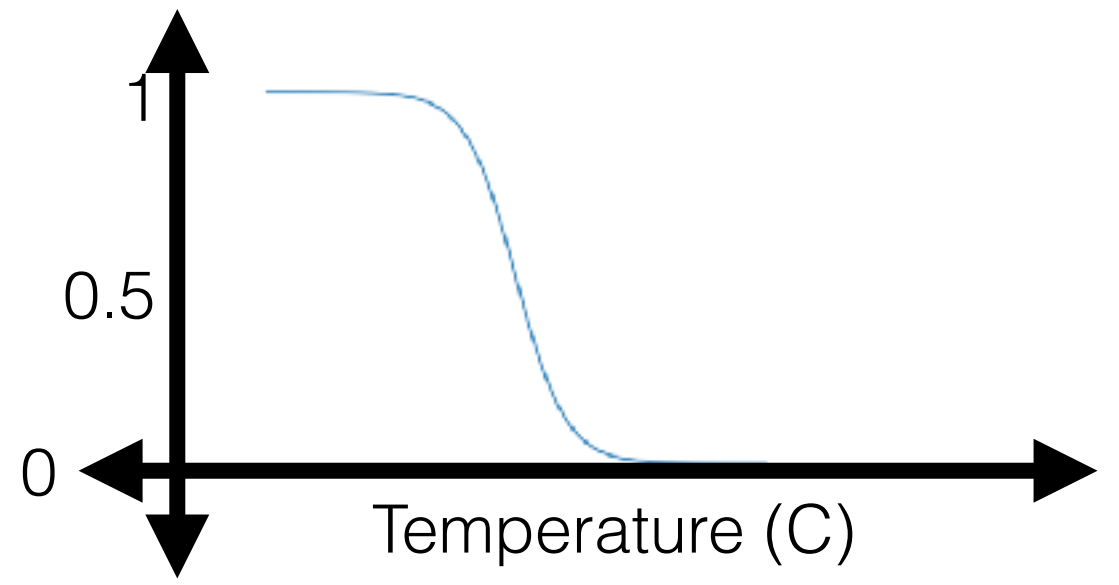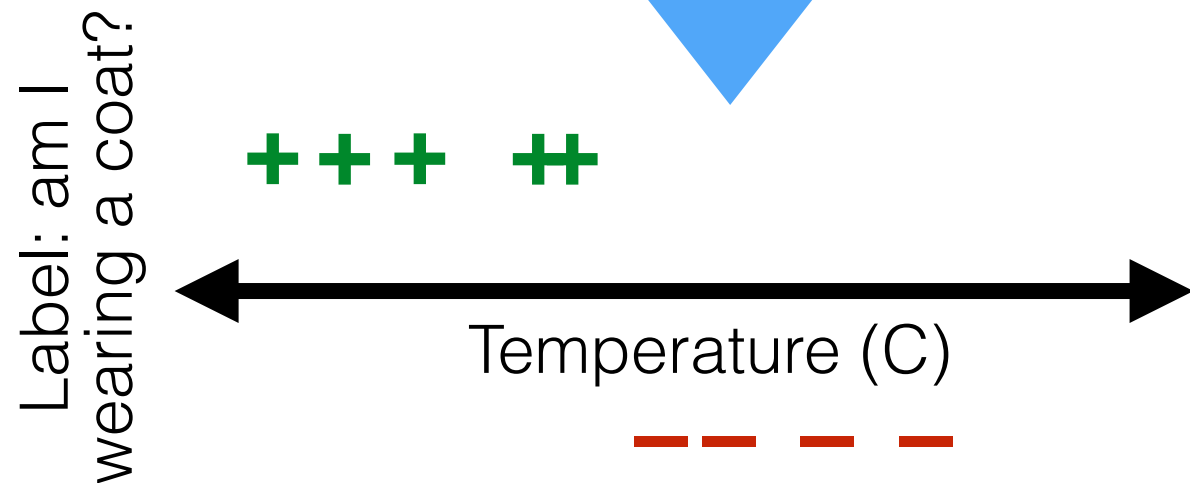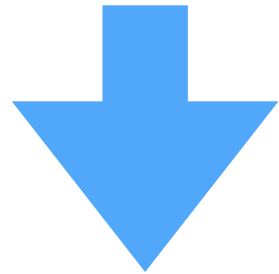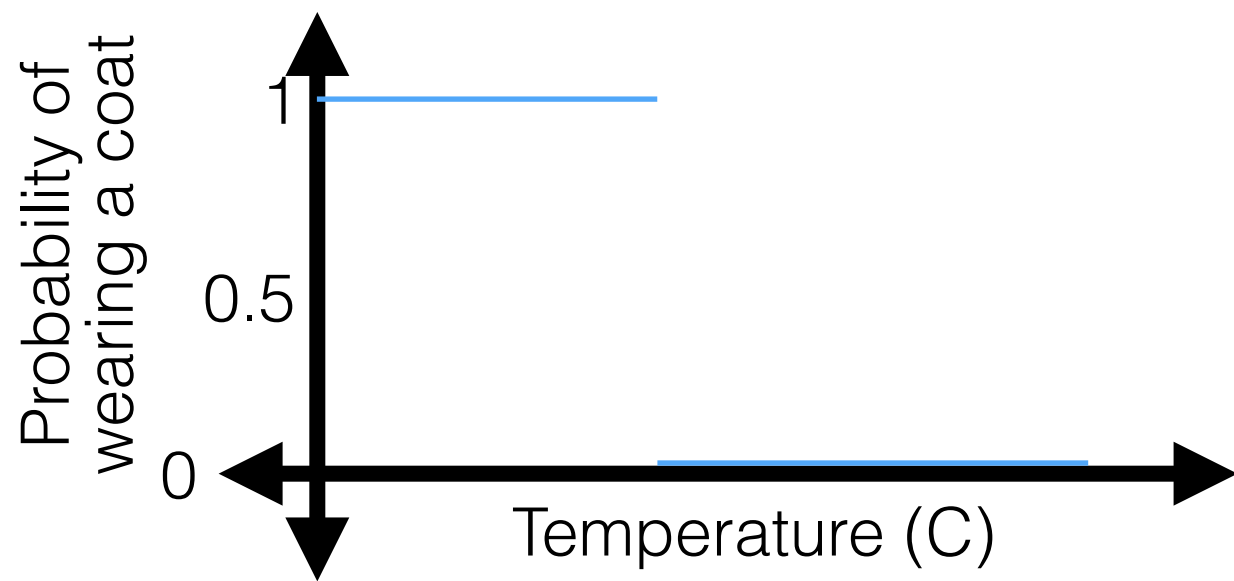
$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

# Capturing uncertainty

Probability of wearing a coat

1

0.5

0

Temperature (C)

1

0.5

0

Temperature (C)

Label: am I wearing a coat?

+ + + ++

Temperature (C)

− − − −

+ + + +++

Temperature (C)

− − − −

- How to make this shape?
  - Sigmoid/logistic function

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

1

0.5

-3    0    3

$z$

# Capturing uncertainty



- How to make this shape?
  - Sigmoid/logistic function
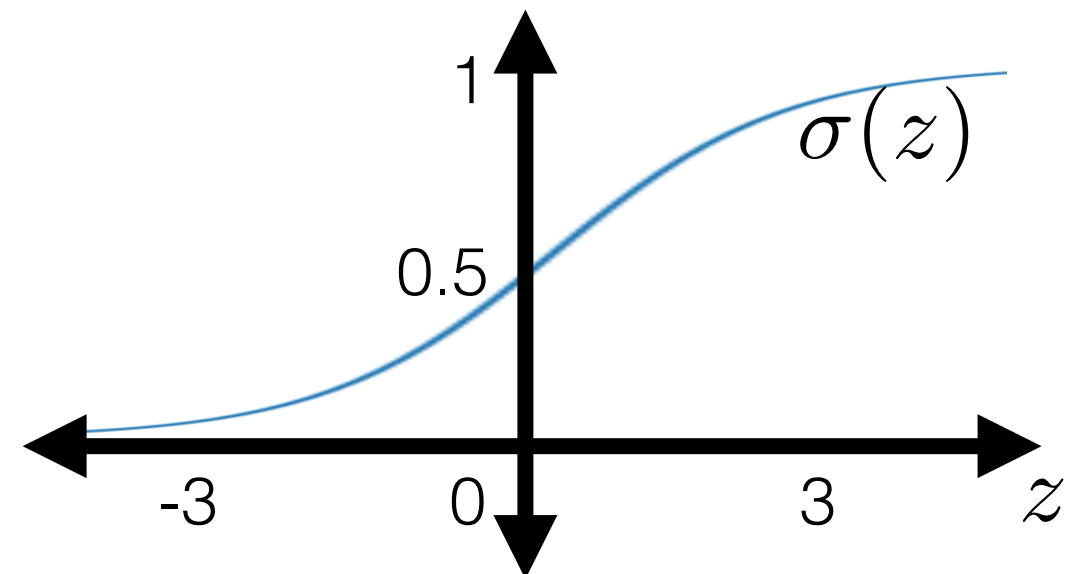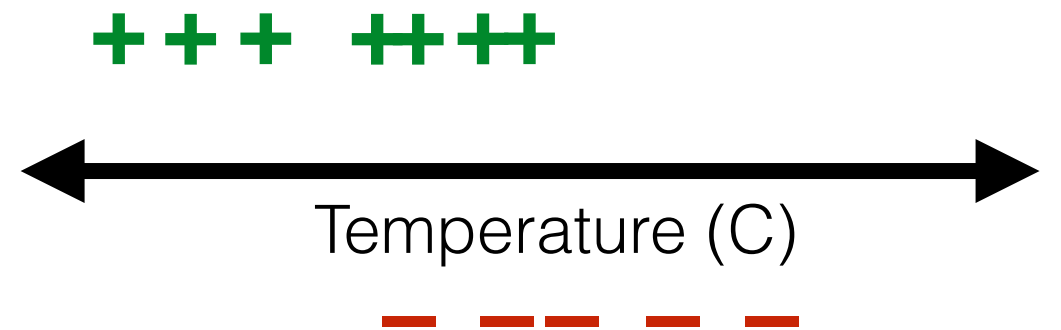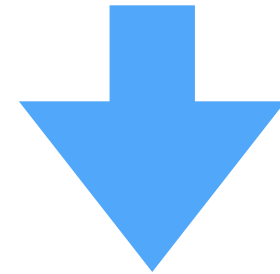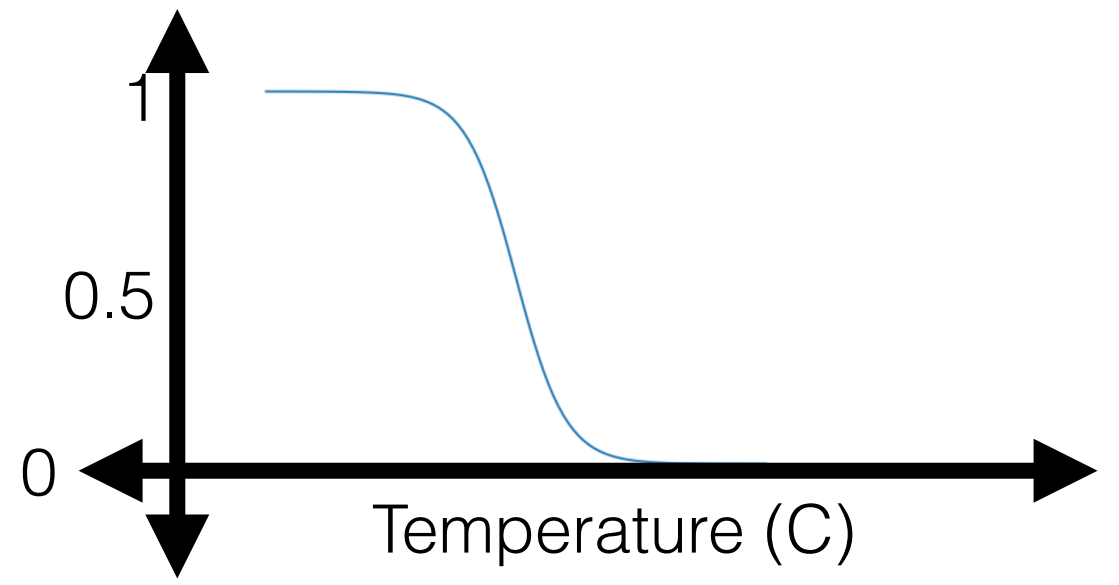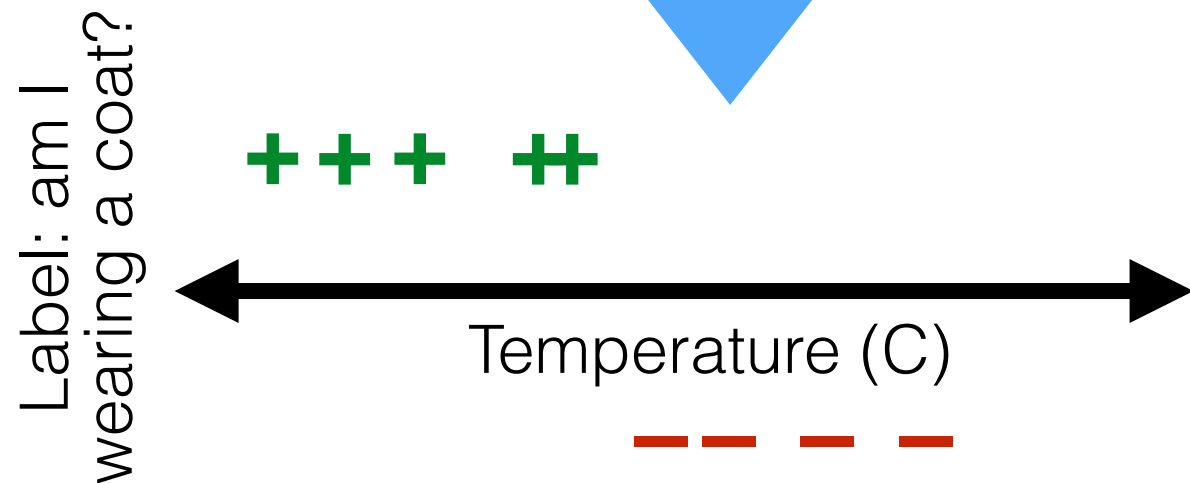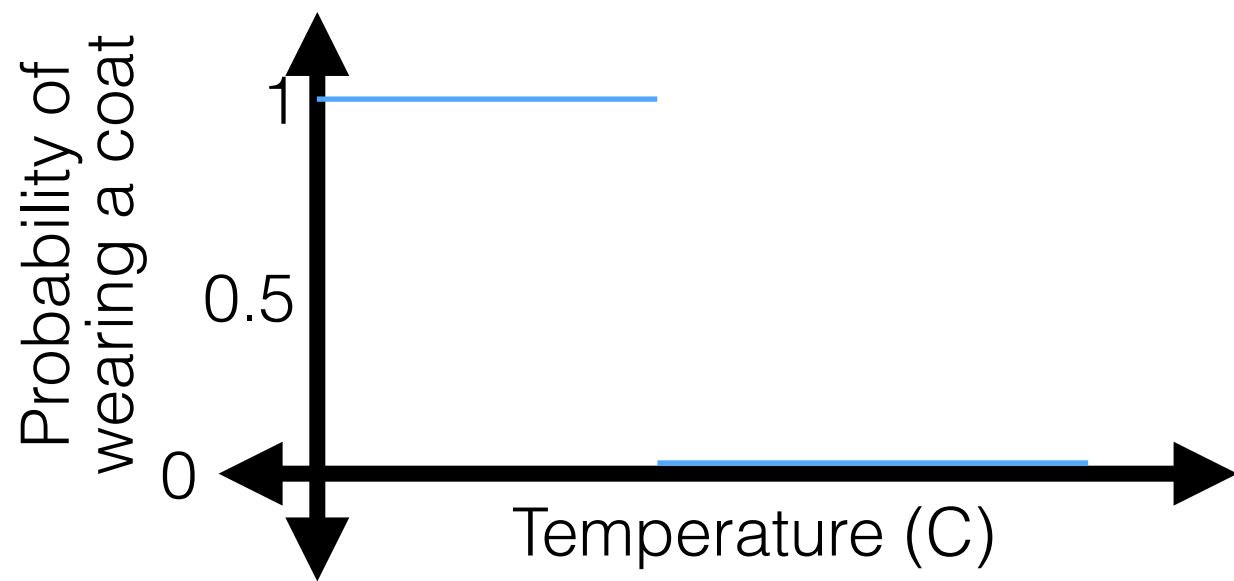
$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

6

# Capturing uncertainty



- How to make this shape?
  - Sigmoid/logistic function

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

6

# Capturing uncertainty
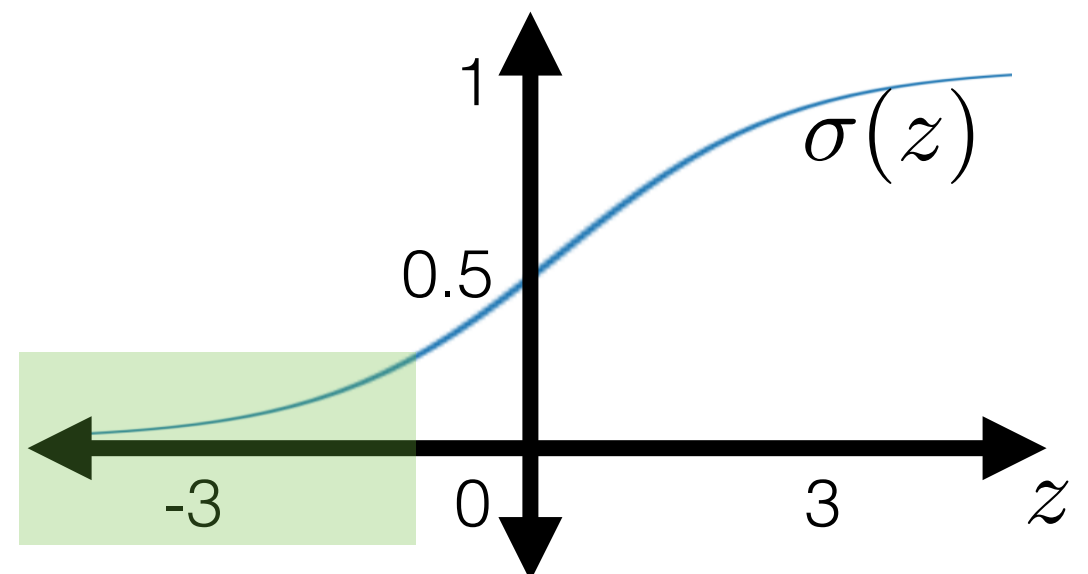
Probability of wearing a coat

1

0.5

0

Temperature (C)

1

0.5

0

Temperature (C)

Label: am I wearing a coat?

+ + + ++

Temperature (C)

− − − −

+ + + ++ +

Temperature (C)

− − − −

- How to make this shape?
  - Sigmoid/logistic function

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

1

$\sigma(z)$

0.5

-3    0    3    $z$

# Capturing uncertainty



- How to make this shape?
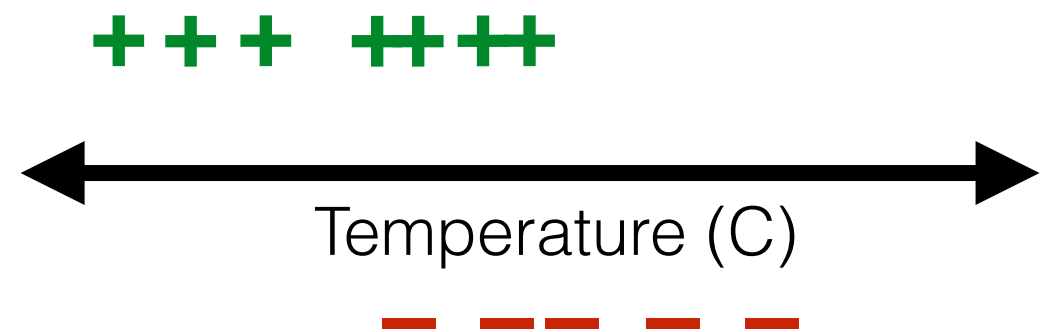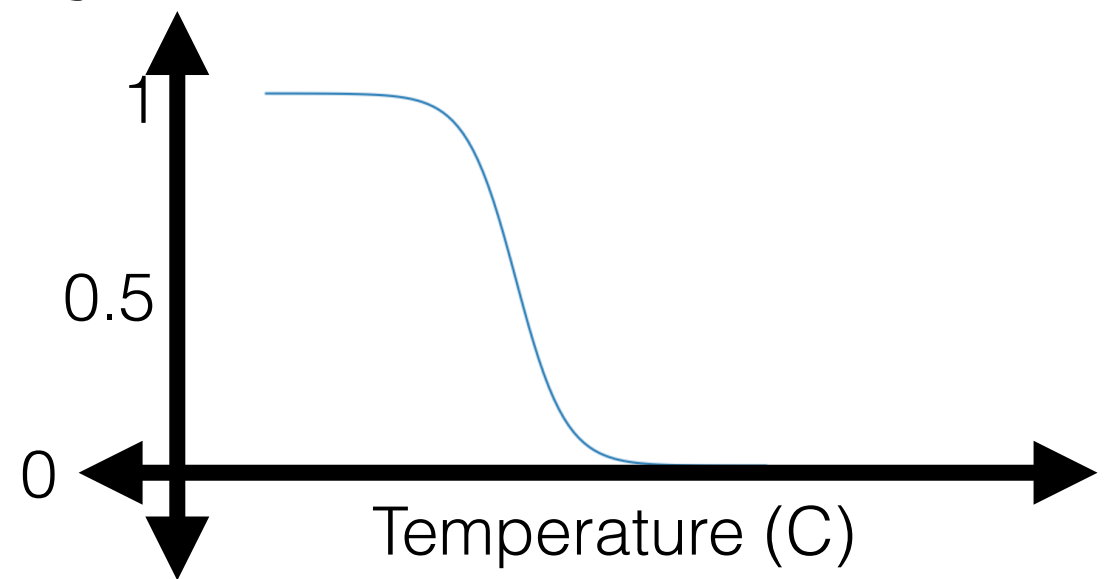  - Sigmoid/logistic function

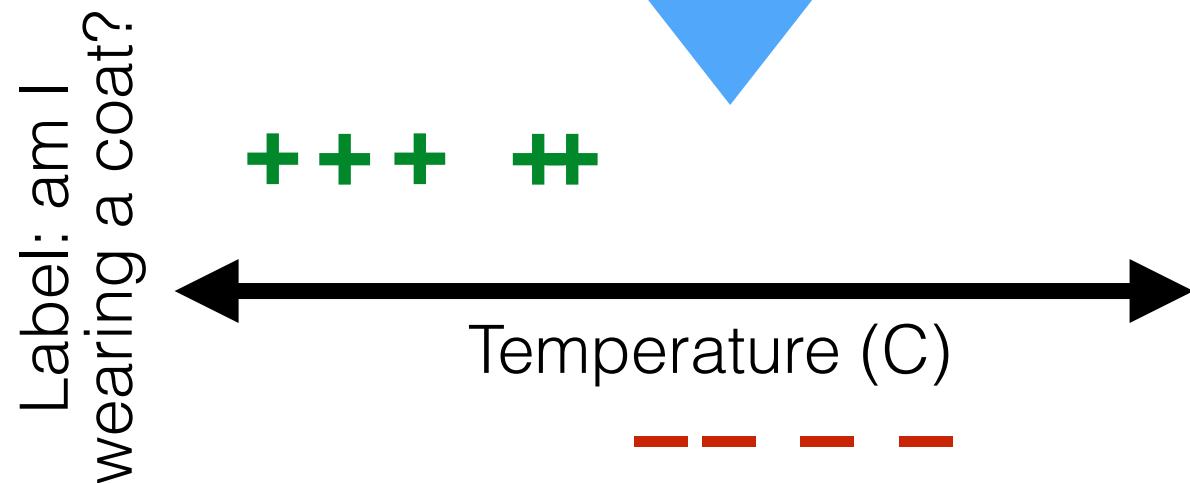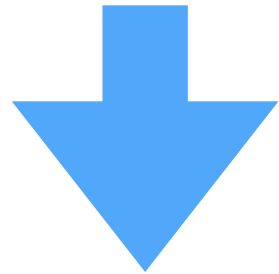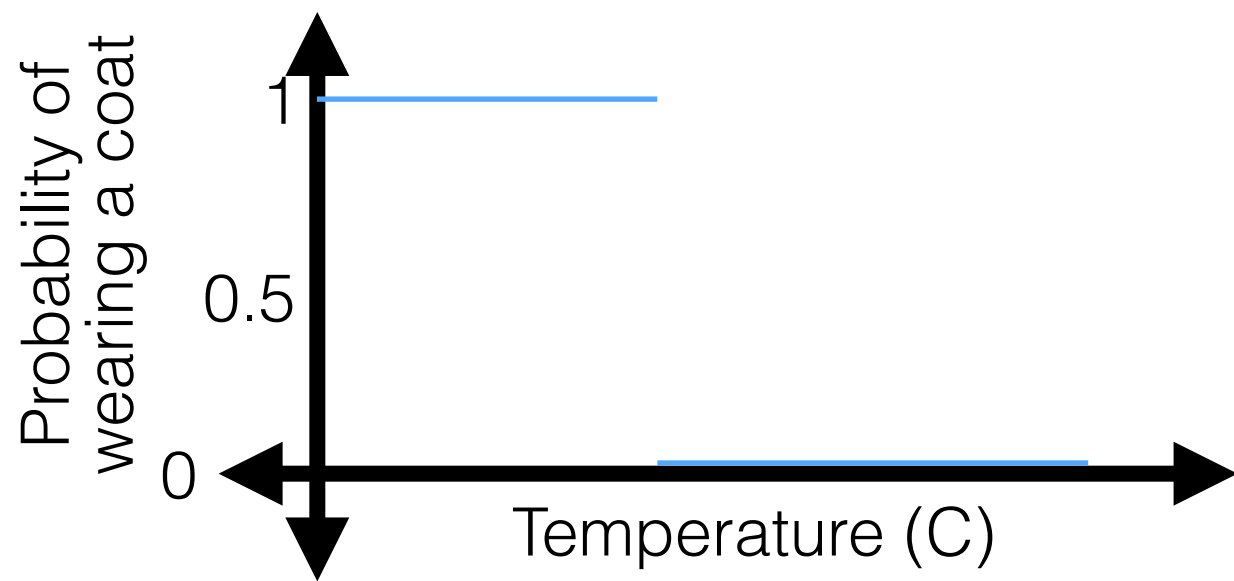$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

# Capturing uncertainty

Probability of wearing a coat

1

0.5

0

Temperature (C)

1

0.5

0

Temperature (C)

Label: am I wearing a coat?

+ + + + + +

Temperature (C)

- - - - - 

+ + + + + + +

Temperature (C)

- - - - -

- How to make this shape?

  - Sigmoid/logistic function

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

$\sigma(z)$

1

0.5

-3     0     3     $z$

# Capturing uncertainty



- How to make this shape?
  - Sigmoid/logistic function
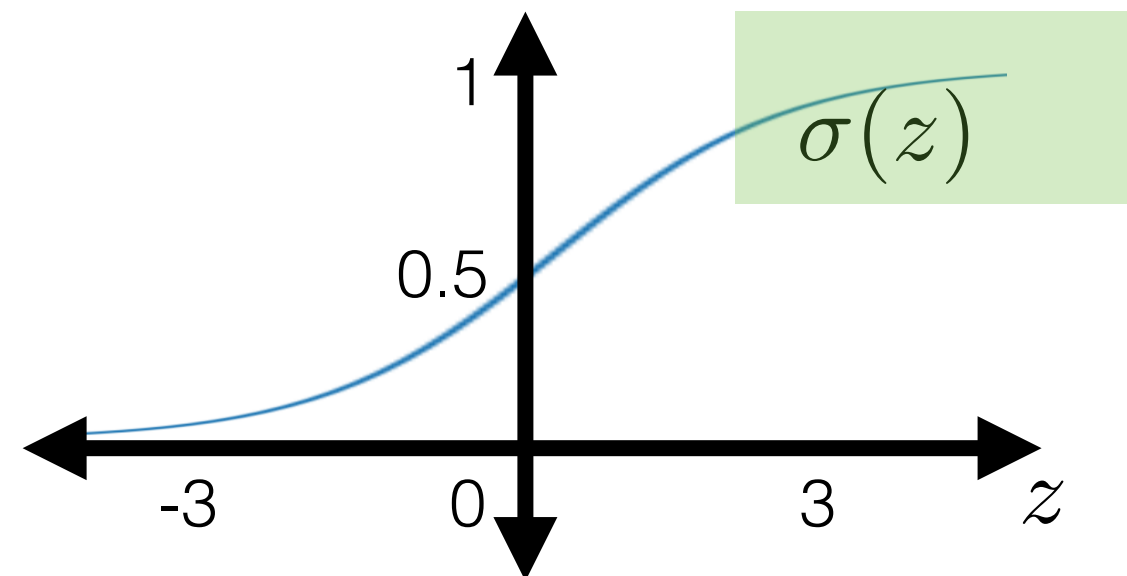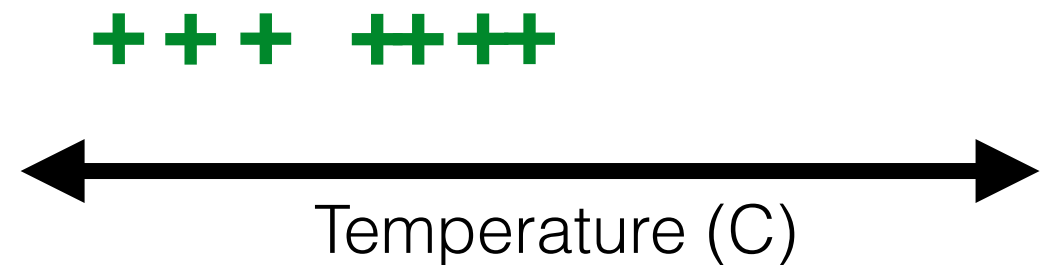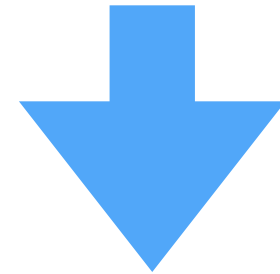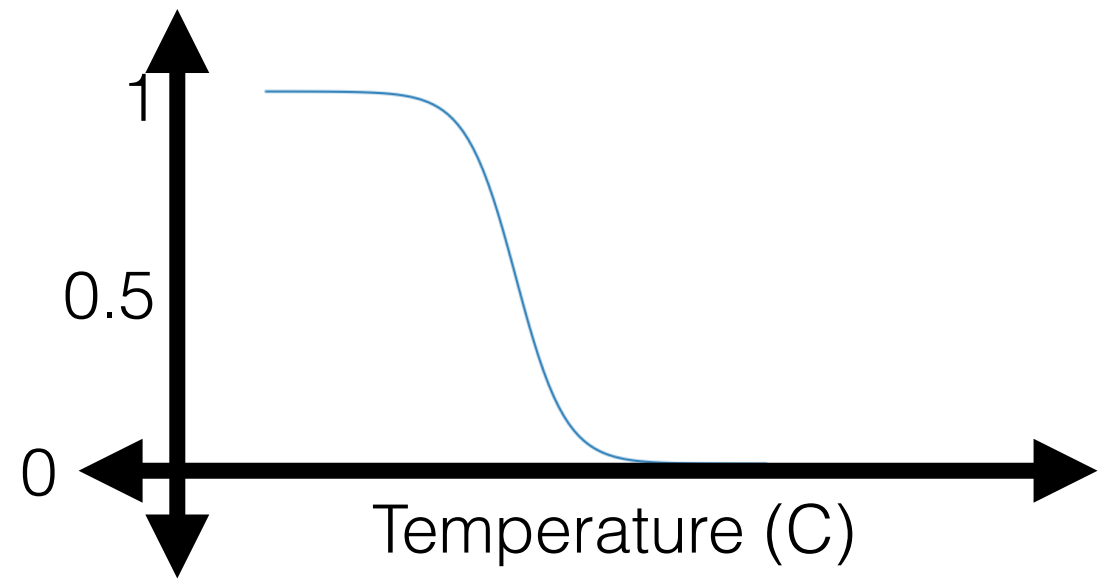
$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

7

# Capturing uncertainty



- How to make this shape?
  - Sigmoid/logistic function

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

# Capturing uncertainty



- How to make this shape?
  - Sigmoid/logistic function

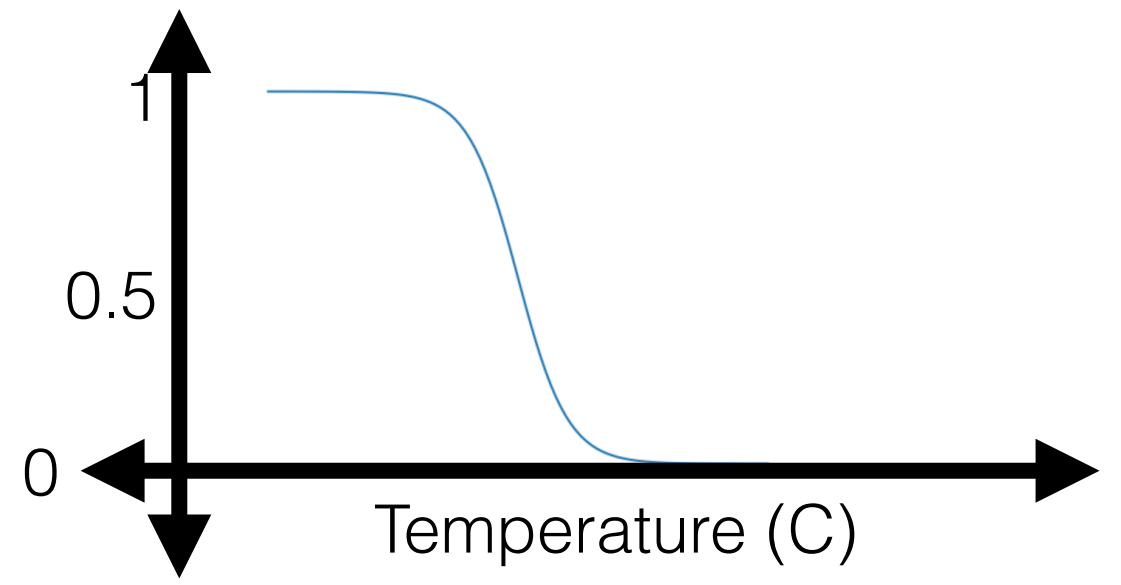$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$
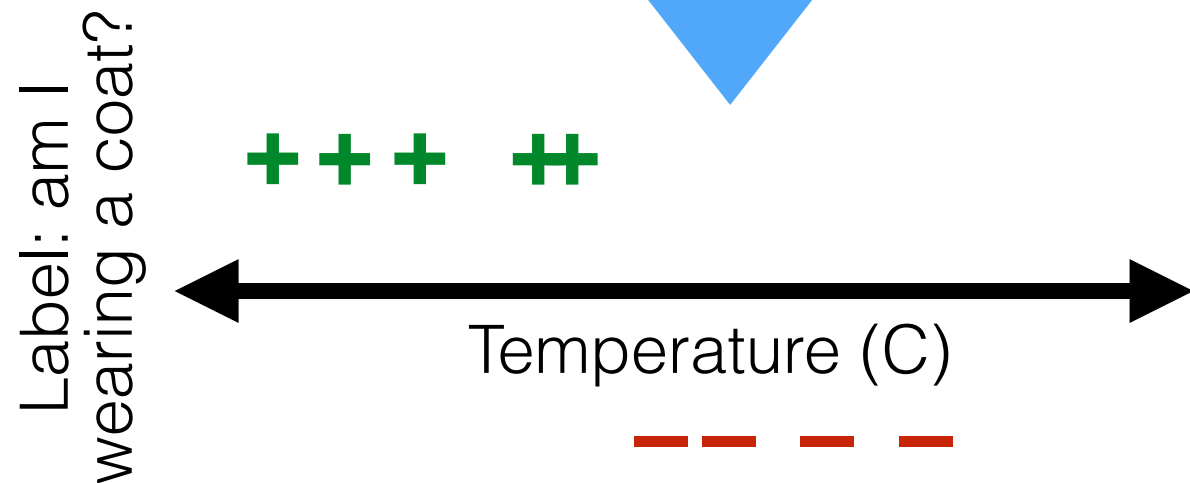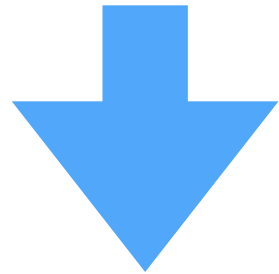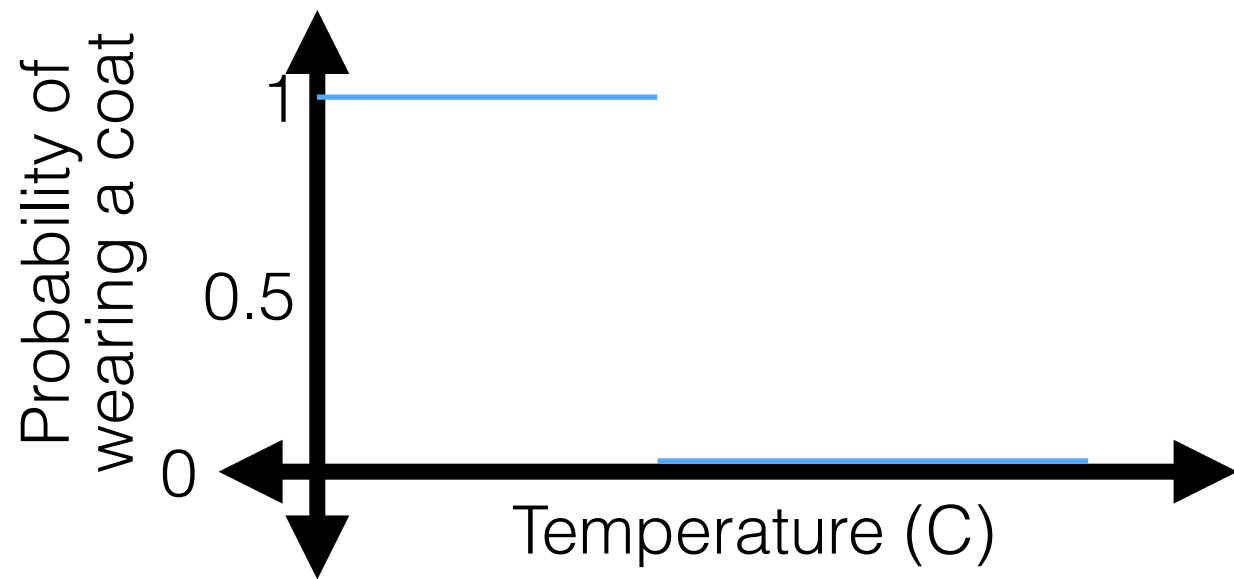
# Capturing uncertainty



- How to make this shape?
  - Sigmoid/logistic function

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

# Capturing uncertainty



- How to make this shape?
  - Sigmoid/logistic function

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

# Capturing uncertainty



- How to make this shape?
  - Sigmoid/logistic function

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

Temperature (C)

+ + + + + +

Temperature (C)

$\sigma(\theta z)$

1

0.5

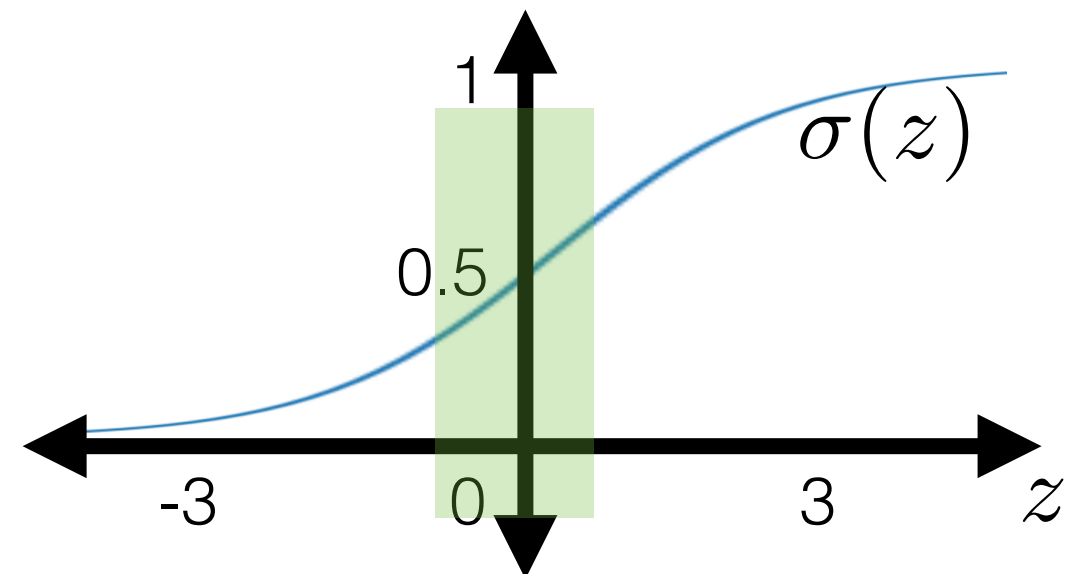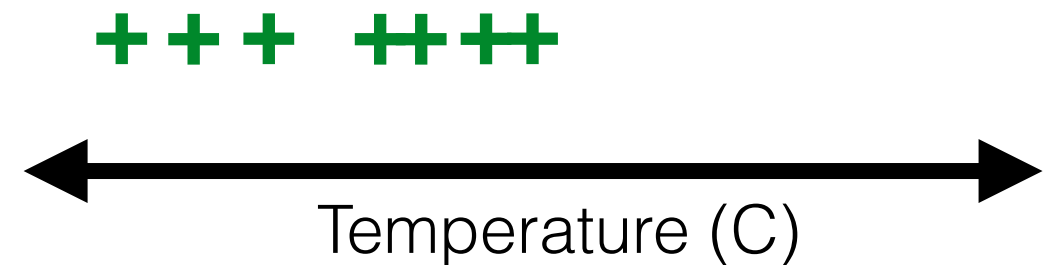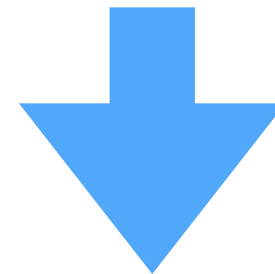-3    0    3    $z$
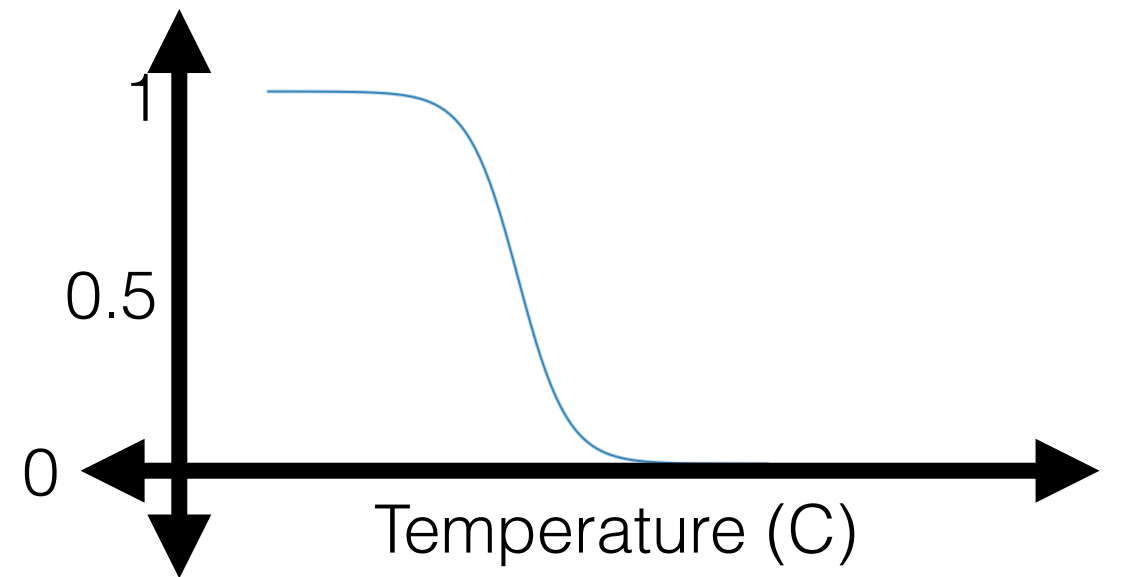
7

# Capturing uncertainty



- How to make this shape?
  - Sigmoid/logistic function

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

7

# Capturing uncertainty



- How to make this shape?
  - Sigmoid/logistic function

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$
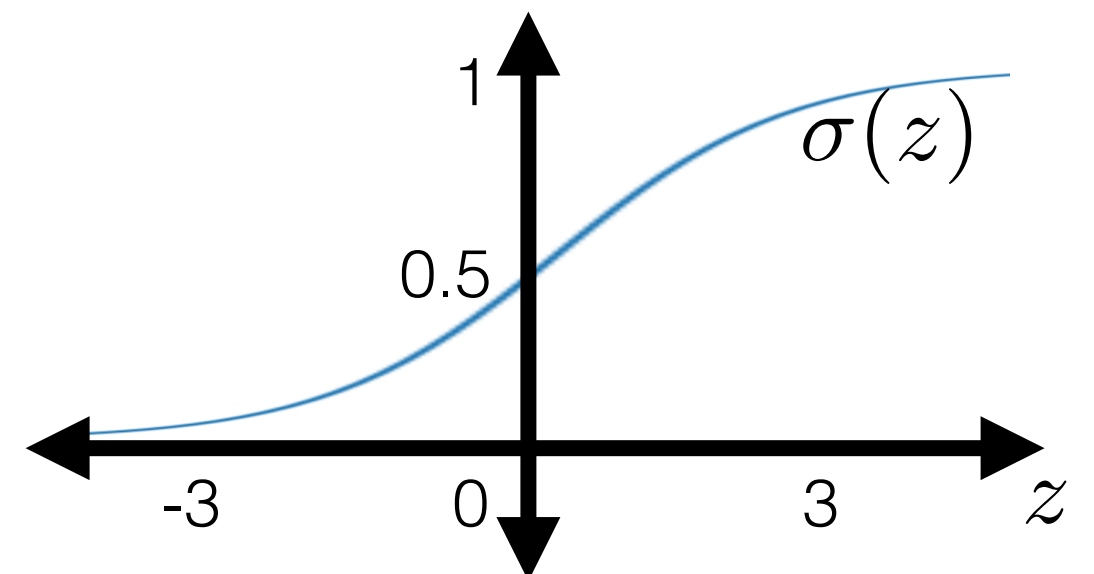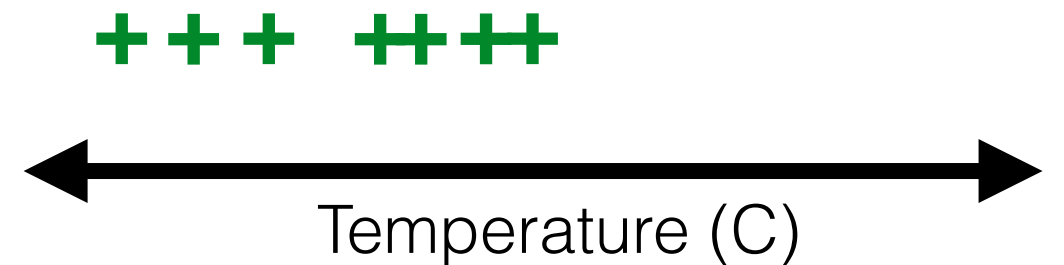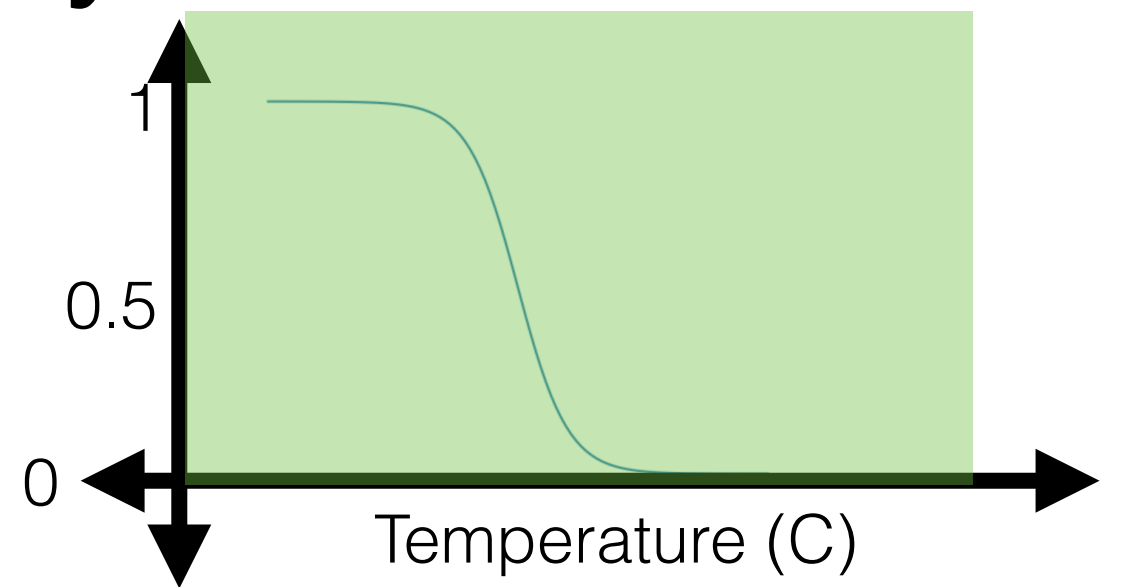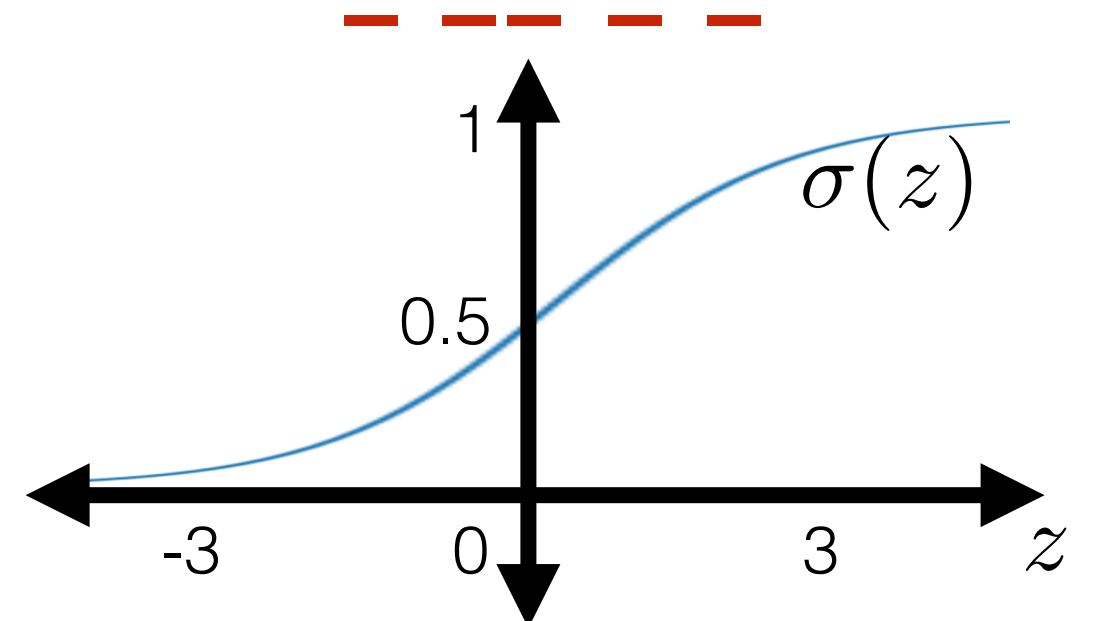
$\sigma(\theta z + \theta_0)$

# Capturing uncertainty



- How to make this shape?
  - Sigmoid/logistic function

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

$\sigma(\theta z + \theta_0)$

# Capturing uncertainty



$g(x)$

1

0.5

0

Temperature (C)

$x$

+ + +  + + +

Temperature (C)

− − − − −

- How to make this shape?
  - Sigmoid/logistic function

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

$\sigma(\theta z + \theta_0)$

1

0.5

-3    0    3    $z$

# Capturing uncertainty

$$g(x) = \sigma(\theta x + \theta_0)$$

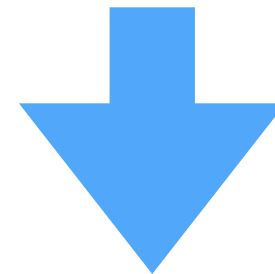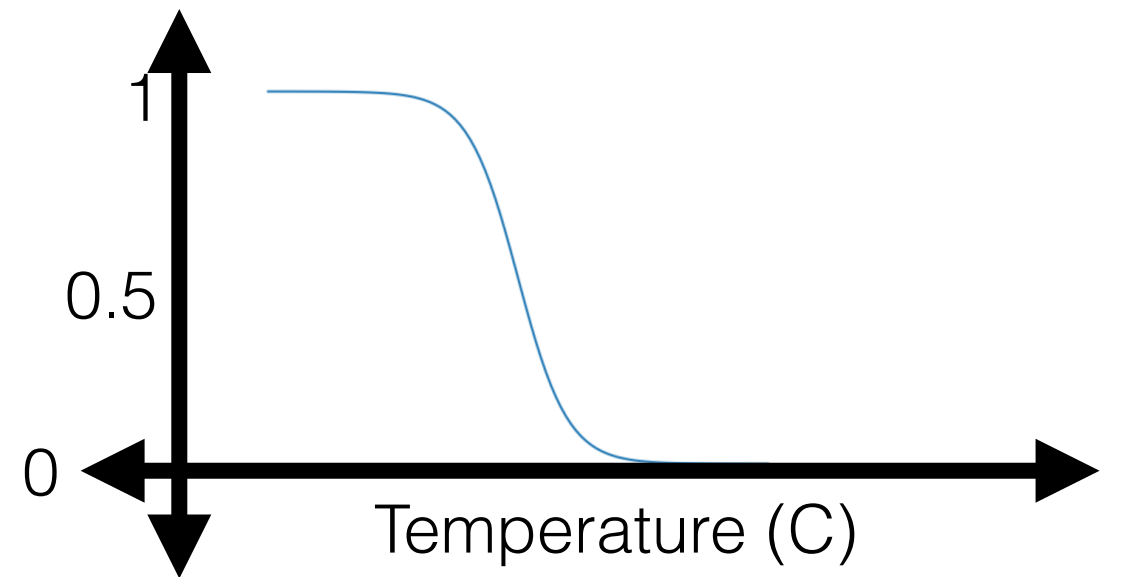$$= \frac{1}{1 + \exp\{-(\theta x + \theta_0)\}}$$
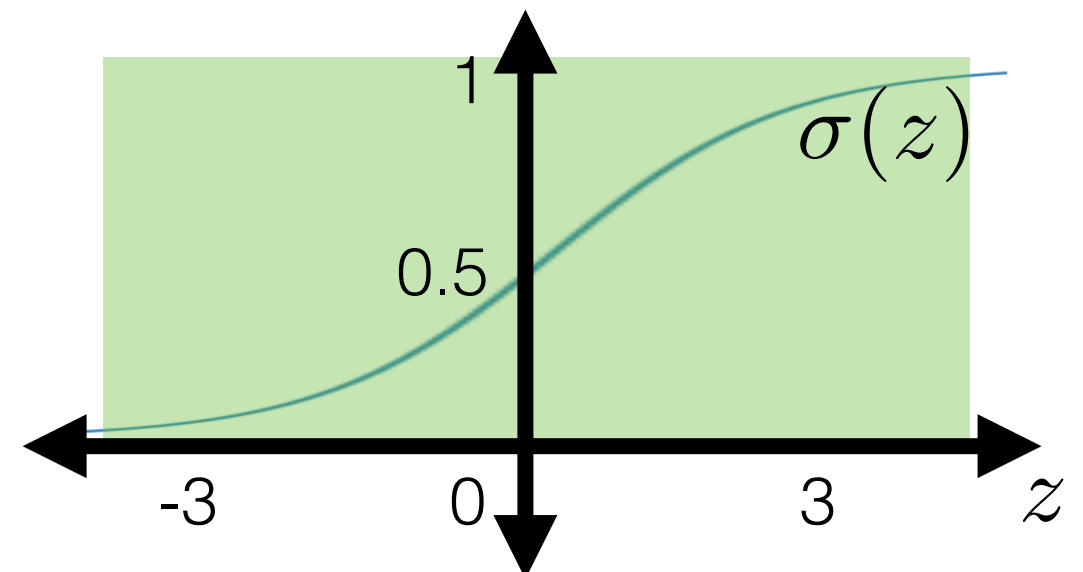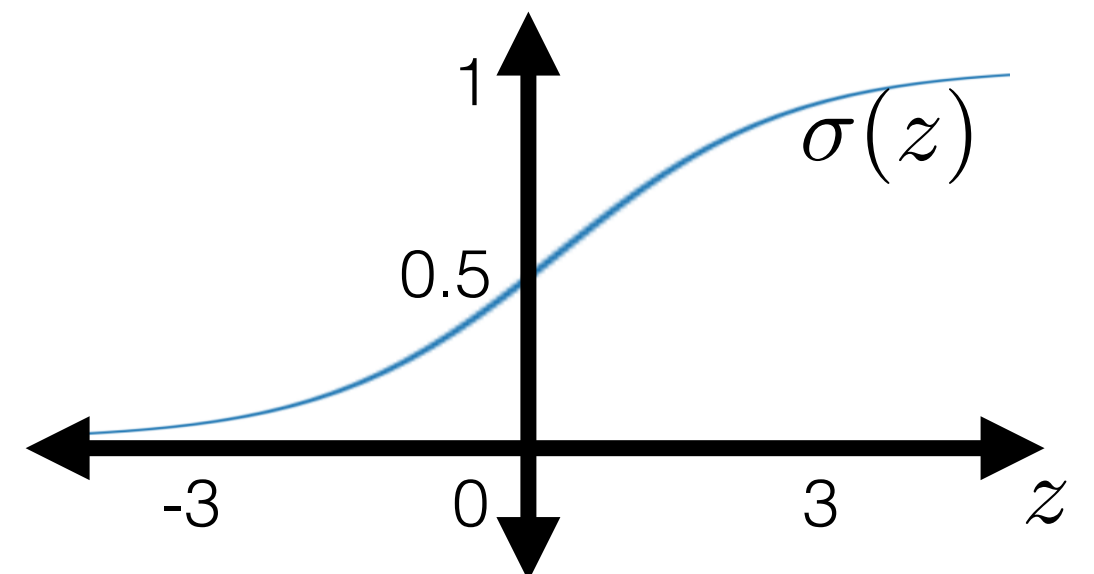


- How to make this shape?
  - Sigmoid/logistic function

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

7

# Capturing uncertainty

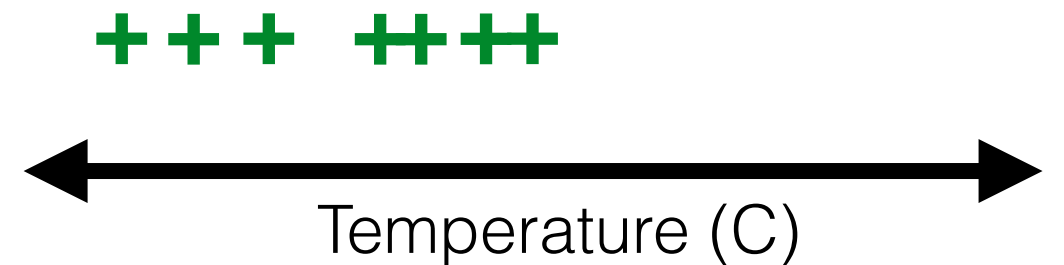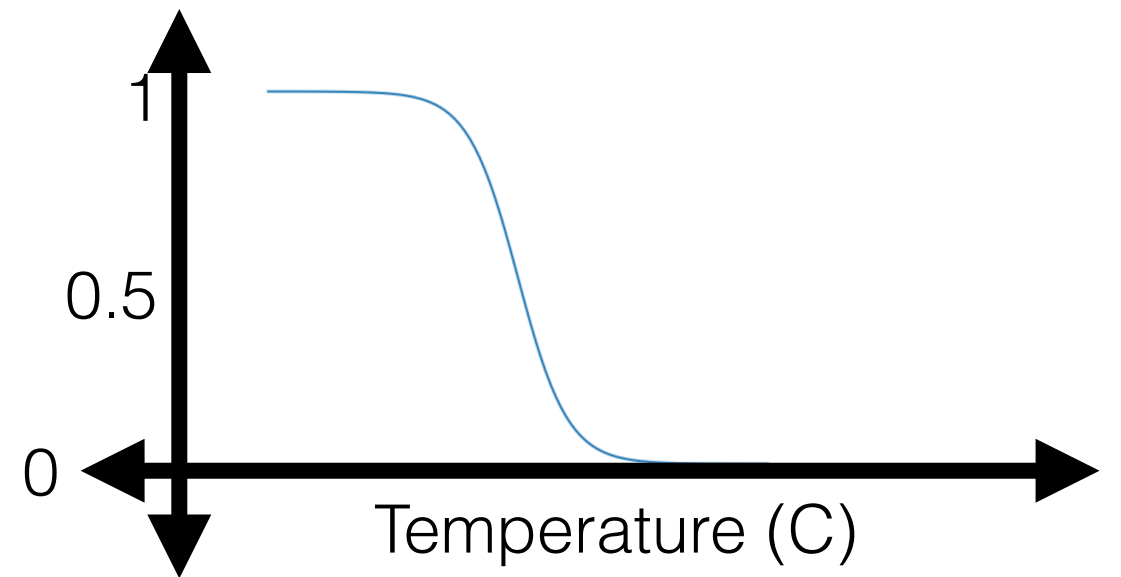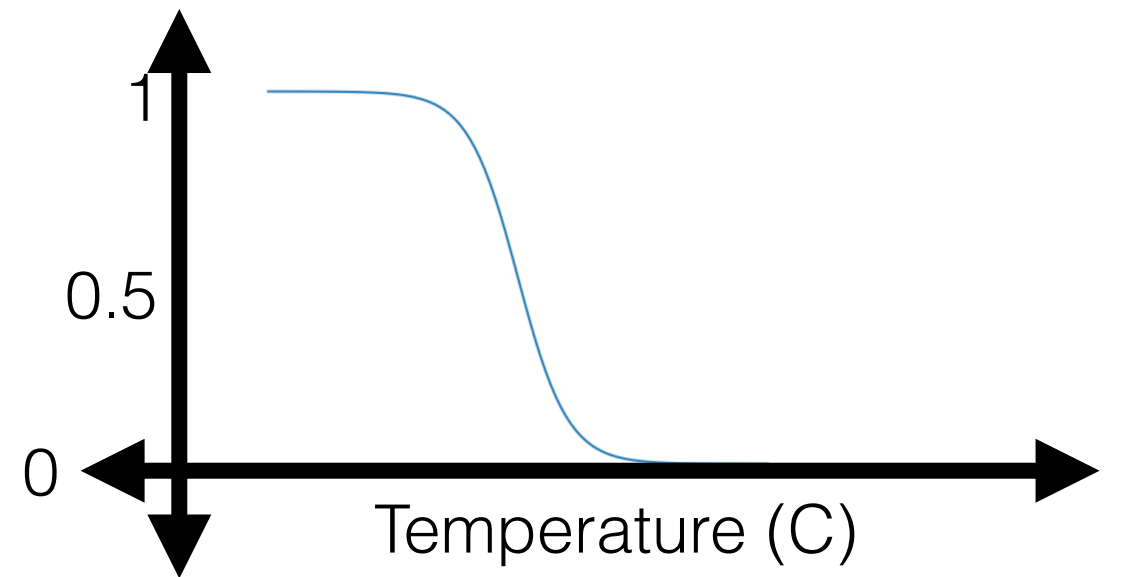$$g(x) = \sigma(\theta^\top x + \theta_0)$$

$$= \frac{1}{1 + \exp\{-(\theta^\top x + \theta_0)\}}$$

$$g(x) = \sigma(\theta x + \theta_0)$$

$$= \frac{1}{1 + \exp\{-(\theta x + \theta_0)\}}$$



$g(x)$

Temperature (C)   $x$



$g(x)$

$x_2$

$x_1$

Temperature (C)

+ + + + + + +

− − − − − 

8

# Capturing uncertainty

$$g(x) = \sigma(\theta x + \theta_0)$$

$$= \frac{1}{1 + \exp\{-(\theta x + \theta_0)\}}$$

$$g(x) = \sigma(\theta^\top x + \theta_0)$$

$$= \frac{1}{1 + \exp\{-(\theta^\top x + \theta_0)\}}$$



8

# Capturing uncertainty

$$g(x) = \sigma(\theta x + \theta_0)$$

$$= \frac{1}{1 + \exp\{-(\theta x + \theta_0)\}}$$

$$g(x) = \sigma(\theta^\top x + \theta_0)$$

$$= \frac{1}{1 + \exp\{-(\theta^\top x + \theta_0)\}}$$



8

# Capturing uncertainty

$$g(x) = \sigma(\theta^\top x + \theta_0)$$

$$= \frac{1}{1 + \exp\{-(\theta^\top x + \theta_0)\}}$$

1 feature:

$$g(x) = \sigma(\theta x + \theta_0)$$

$$= \frac{1}{1 + \exp\{-(\theta x + \theta_0)\}}$$



8

# Linear logistic classification

- What's an appropriate loss for this guess?

# Linear logistic classification

- What's an appropriate loss for this guess?

# Linear logistic classification

- What's an appropriate loss for this guess?

# Linear logistic classification

- What's an appropriate loss for this guess?

# Linear logistic classification

- What's an appropriate loss for this guess?

Probability(data)

$$= \prod_{i=1}^{n} \text{Probability(data point } i)$$

# Linear logistic classification

- What's an appropriate loss for this guess?

Probability(data)

$$= \prod_{i=1}^{n} \text{Probability(data point } i)$$

$$[\text{Let } g^{(i)} = \sigma(\theta^\top x^{(i)} + \theta_0)]$$

$$= \prod_{i=1}^{n} \begin{cases} g^{(i)} \text{ if } y^{(i)} = +1 \\ (1 - g^{(i)}) \text{ else} \end{cases}$$

$$= \prod_{i=1}^{n} (g^{(i)})^{\mathbf{1}\{y^{(i)}=+1\}} (1 - g^{(i)})^{\mathbf{1}\{y^{(i)}\neq+1\}}$$



9

# Linear logistic classification

- What's an appropriate loss for this guess?

Probability(data)

$$= \prod_{i=1}^{n} \text{Probability}(\text{data point } i)$$

$$\quad [\text{Let } g^{(i)} = \sigma(\theta^\top x^{(i)} + \theta_0)]$$

$$= \prod_{i=1}^{n} \begin{cases} g^{(i)} \text{ if } y^{(i)} = +1 \\ (1 - g^{(i)}) \text{ else} \end{cases}$$

$$= \prod_{i=1}^{n} (g^{(i)})^{\mathbf{1}\{y^{(i)}=+1\}} (1 - g^{(i)})^{\mathbf{1}\{y^{(i)}\neq+1\}}$$

Temperature (C)

Temperature (C)

$+ + +$ $+ + +$

$- - - - -$

9

# Linear logistic classification

- What's an appropriate loss for this guess?

Probability(data)

$$= \prod_{i=1}^{n} \text{Probability(data point } i)$$

$$[\text{Let } g^{(i)} = \sigma(\theta^{\top} x^{(i)} + \theta_0)]$$

$$= \prod_{i=1}^{n} \begin{cases} g^{(i)} \text{ if } y^{(i)} = +1 \\ (1 - g^{(i)}) \text{ else} \end{cases}$$

$$= \prod_{i=1}^{n} (g^{(i)})^{\mathbf{1}\{y^{(i)}=+1\}} (1 - g^{(i)})^{\mathbf{1}\{y^{(i)}\neq+1\}}$$

Temperature (C)

+ + +   + + +

Temperature (C)

− −− − −

9

# Linear logistic classification

- What's an appropriate loss for this guess?

Probability(data)

$$= \prod_{i=1}^{n} \text{Probability}(\text{data point } i)$$

$$[\text{Let } g^{(i)} = \sigma(\theta^{\top} x^{(i)} + \theta_0)]$$

$$= \prod_{i=1}^{n} \begin{cases} g^{(i)} \text{ if } y^{(i)} = +1 \\ (1 - g^{(i)}) \text{ else} \end{cases}$$

$$= \prod_{i=1}^{n} (g^{(i)})^{\mathbf{1}\{y^{(i)} = +1\}} (1 - g^{(i)})^{\mathbf{1}\{y^{(i)} \neq +1\}}$$

Loss(data) = -log probability(data)

$$= \sum_{i=1}^{n} - \left( \mathbf{1}\{y^{(i)} = +1\} \log g^{(i)} + \mathbf{1}\{y^{(i)} \neq +1\} \log(1 - g^{(i)}) \right)$$



Temperature (C)

+ + +   + + +

Temperature (C)

− −− − −

# Linear logistic classification

- What's an appropriate loss for this guess?

Probability(data)

$$= \prod_{i=1}^{n} \text{Probability(data point } i)$$

$$\left[\text{Let } g^{(i)} = \sigma(\theta^\top x^{(i)} + \theta_0)\right]$$

$$= \prod_{i=1}^{n} \begin{cases} g^{(i)} \text{ if } y^{(i)} = +1 \\ (1 - g^{(i)}) \text{ else} \end{cases}$$

$$= \prod_{i=1}^{n} (g^{(i)})^{\mathbf{1}\{y^{(i)}=+1\}} (1 - g^{(i)})^{\mathbf{1}\{y^{(i)} \neq +1\}}$$

Loss(data) = $-$log probability(data)

$$= \sum_{i=1}^{n} - \left( \mathbf{1}\{y^{(i)} = +1\} \log g^{(i)} + \mathbf{1}\{y^{(i)} \neq +1\} \log(1 - g^{(i)}) \right)$$



Temperature (C)

Temperature (C)

9

# Linear logistic classification
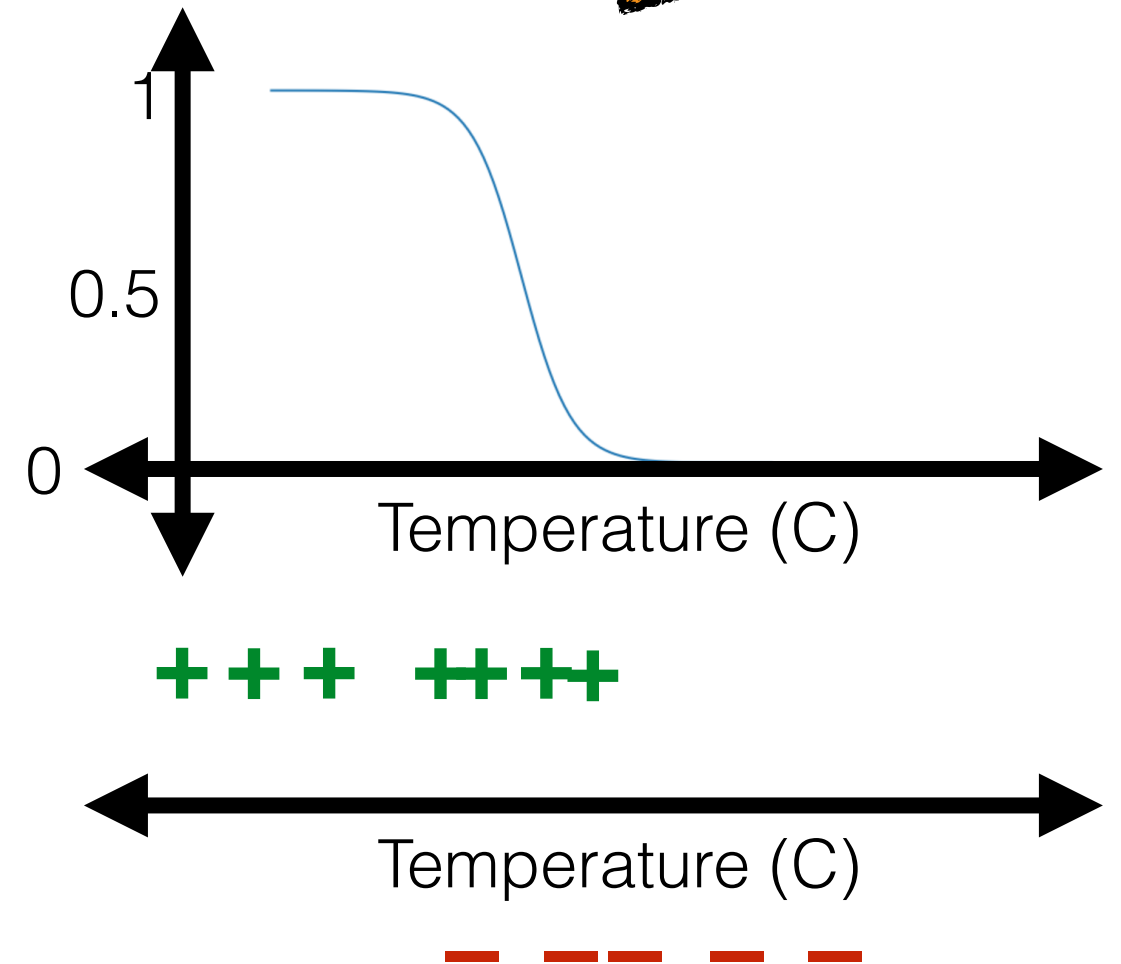
- What's an appropriate loss for this guess?

Probability(data)

$$= \prod_{i=1}^{n} \text{Probability(data point } i)$$

$$[\text{Let } g^{(i)} = \sigma(\theta^{\top} x^{(i)} + \theta_0)]$$

$$= \prod_{i=1}^{n} \begin{cases} g^{(i)} \text{ if } y^{(i)} = +1 \\ (1 - g^{(i)}) \text{ else} \end{cases}$$

$$= \prod_{i=1}^{n} (g^{(i)})^{\mathbf{1}\{y^{(i)}=+1\}} (1 - g^{(i)})^{\mathbf{1}\{y^{(i)}\neq+1\}}$$

1

0.5

0

Temperature (C)

+ + +  + + +

Temperature (C)

− −− − −

Loss(data) =      - log probability(data)

$$= \sum_{i=1}^{n} - \left( \mathbf{1}\{y^{(i)} = +1\} \log g^{(i)} + \mathbf{1}\{y^{(i)} \neq +1\} \log(1 - g^{(i)}) \right)$$

9

# Linear logistic classification
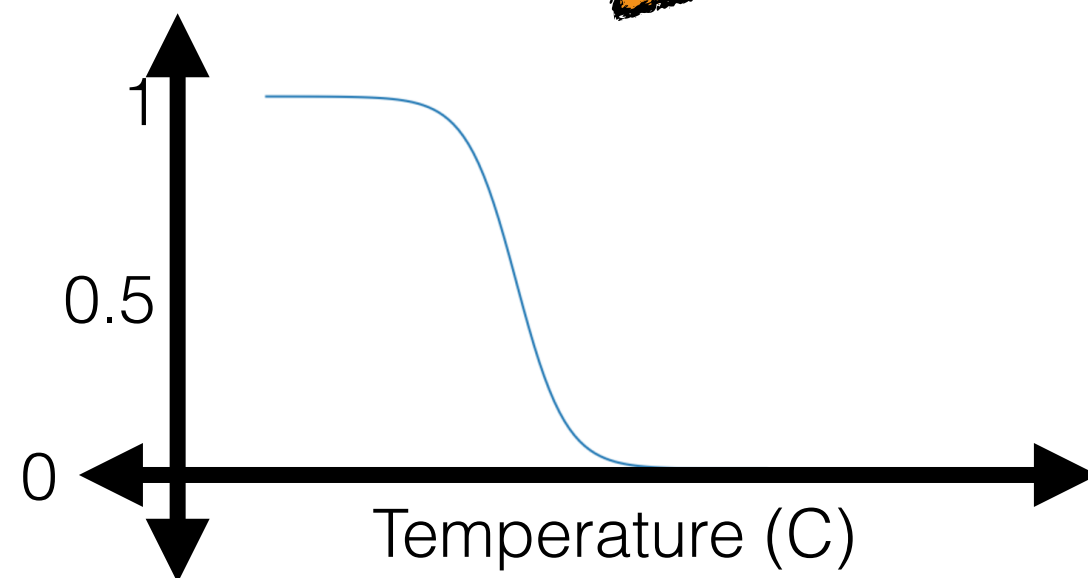
- What's an appropriate loss for this guess?

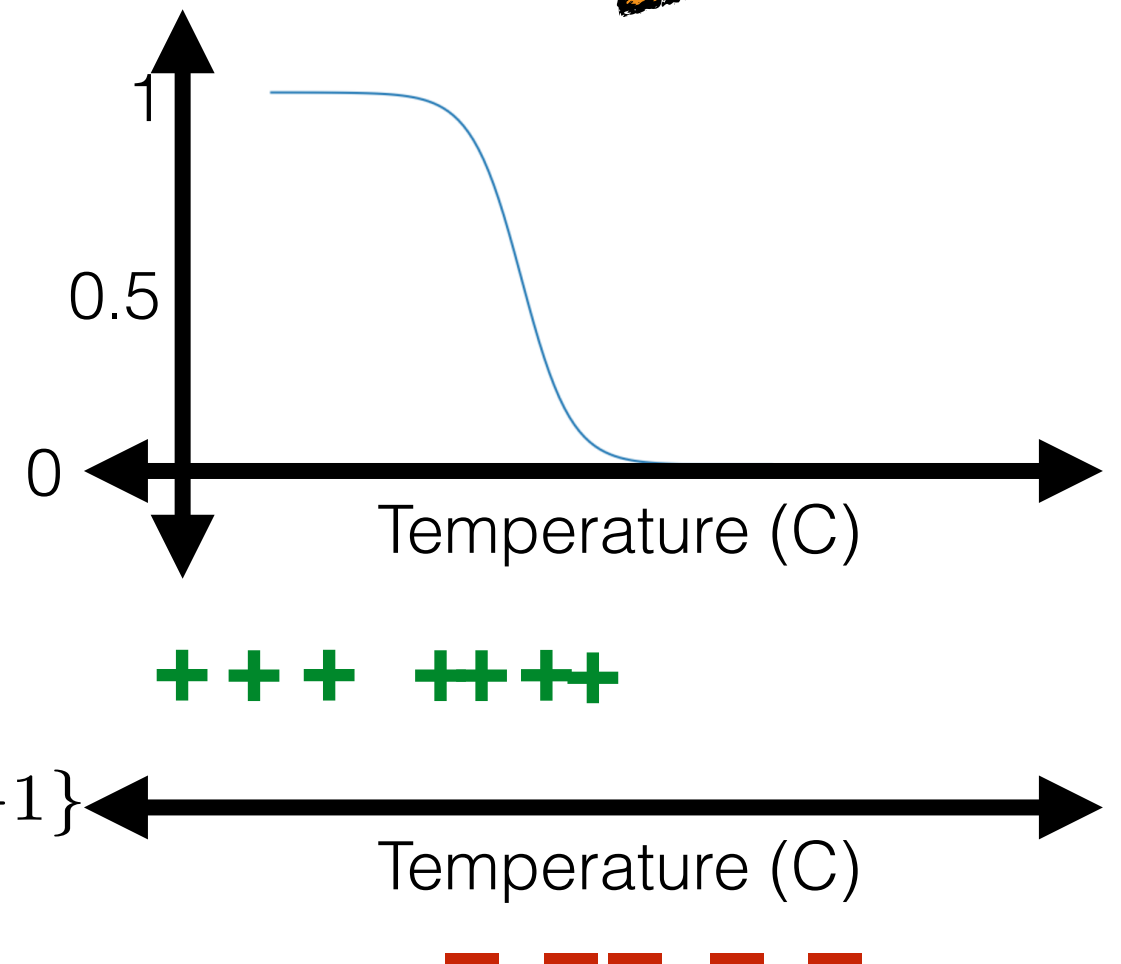Probability(data)

$$= \prod_{i=1}^{n} \text{Probability}(\text{data point } i)$$
$$[\text{Let } g^{(i)} = \sigma(\theta^{\top} x^{(i)} + \theta_0)]$$

$$= \prod_{i=1}^{n} \begin{cases} g^{(i)} \text{ if } y^{(i)} = +1 \\ (1 - g^{(i)}) \text{ else} \end{cases}$$

$$= \prod_{i=1}^{n} (g^{(i)})^{\mathbf{1}\{y^{(i)}=+1\}} (1 - g^{(i)})^{\mathbf{1}\{y^{(i)}\neq+1\}}$$

Loss(data) = $-\log$ probability(data)

$$= \sum_{i=1}^{n} - \left( \mathbf{1}\{y^{(i)} = +1\} \log g^{(i)} + \mathbf{1}\{y^{(i)} \neq +1\} \log(1 - g^{(i)}) \right.$$

1

0.5

0

Temperature (C)

+ + +   + + +

Temperature (C)

– –– – –

9

# Linear logistic classification

- What's an appropriate loss for this guess?

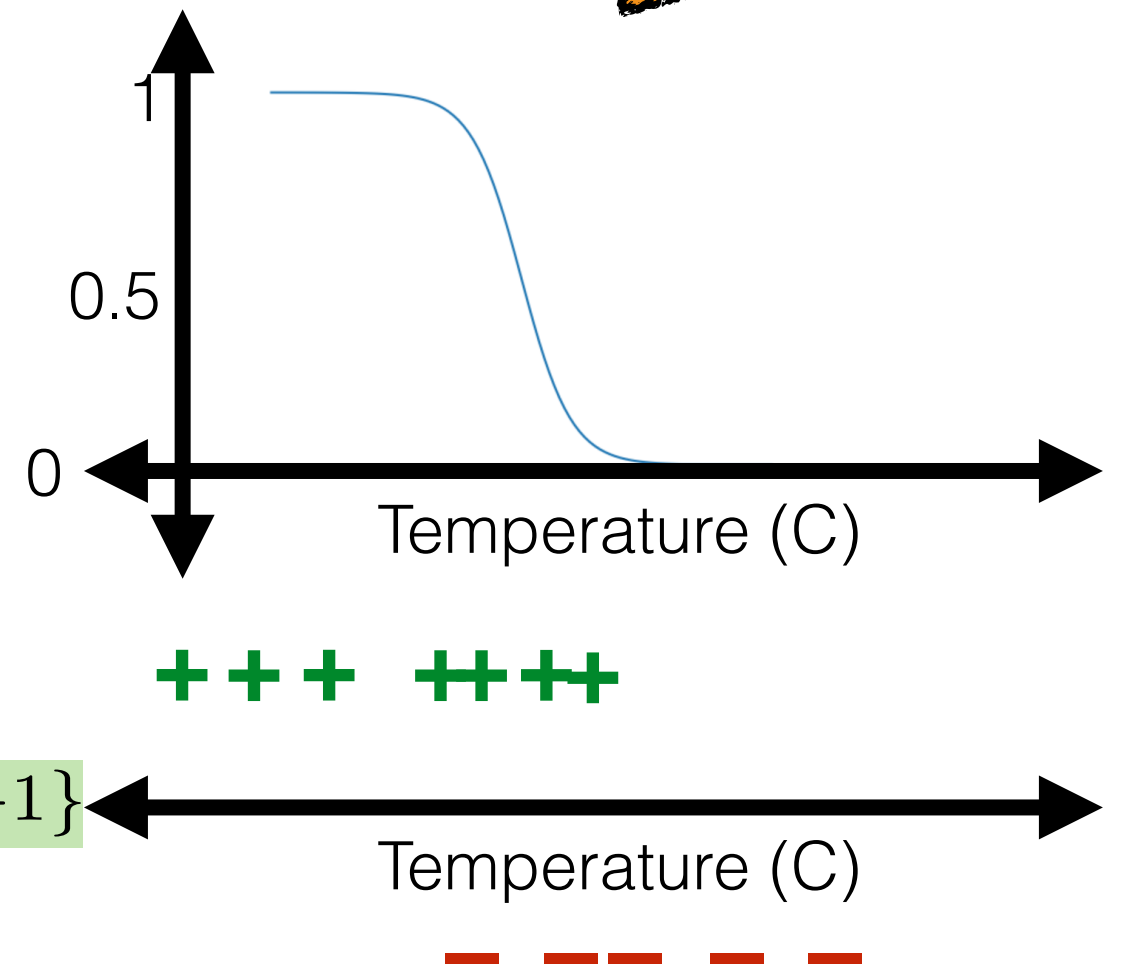Probability(data)

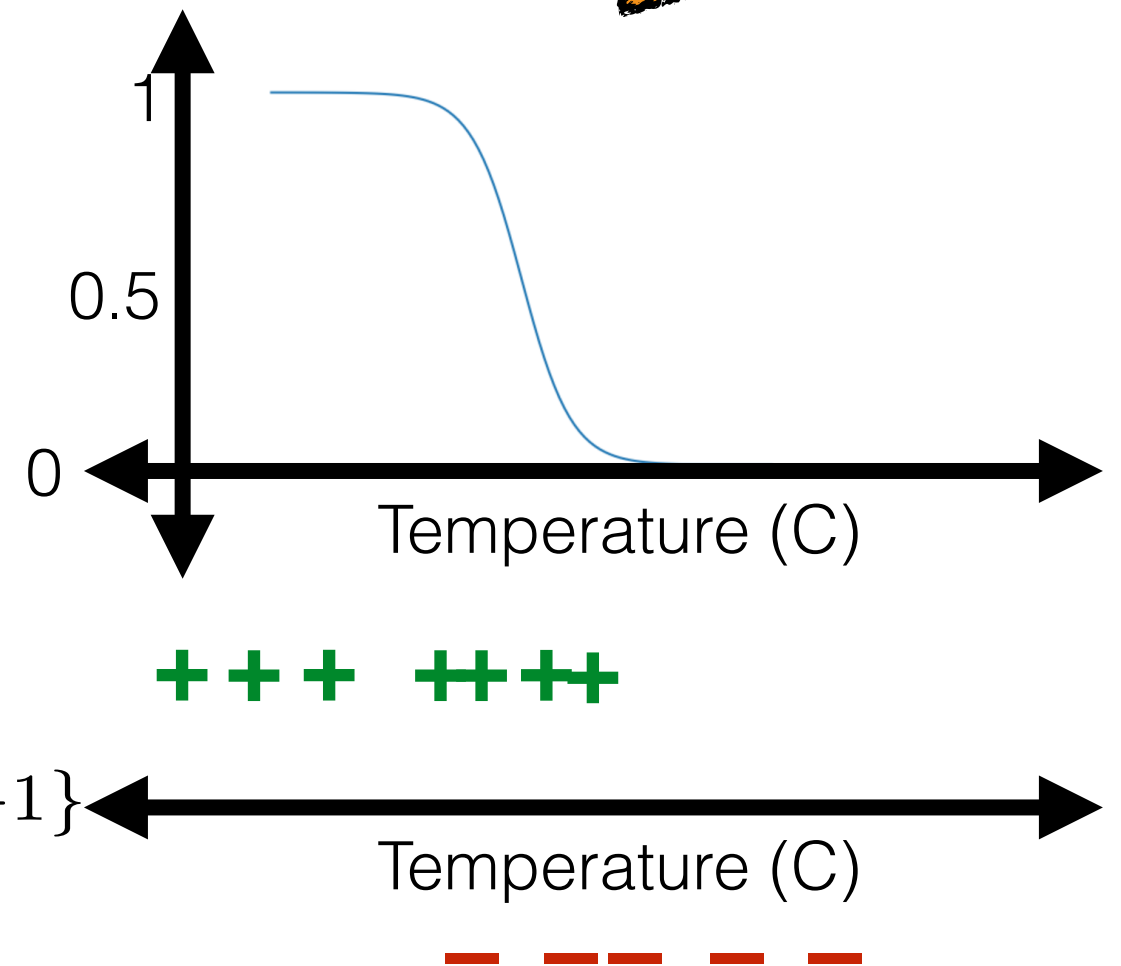$$= \prod_{i=1}^{n} \text{Probability}(\text{data point } i)$$

$$[\text{Let } g^{(i)} = \sigma(\theta^\top x^{(i)} + \theta_0)]$$

$$= \prod_{i=1}^{n} \begin{cases} g^{(i)} & \text{if } y^{(i)} = +1 \\ (1 - g^{(i)}) & \text{else} \end{cases}$$

$$= \prod_{i=1}^{n} (g^{(i)})^{\mathbf{1}\{y^{(i)}=+1\}} (1 - g^{(i)})^{\mathbf{1}\{y^{(i)} \neq +1\}}$$

Loss(data) = $-$log probability(data)

$$= \sum_{i=1}^{n} -\left( \mathbf{1}\{y^{(i)} = +1\} \log g^{(i)} + \mathbf{1}\{y^{(i)} \neq +1\} \log(1 - g^{(i)}) \right)$$

9

# Linear logistic classification
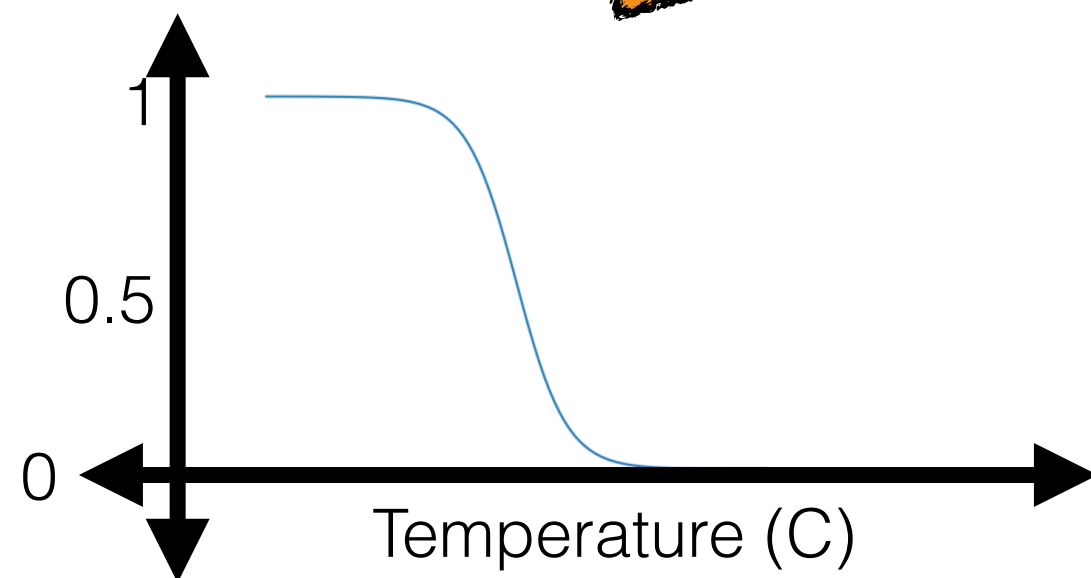
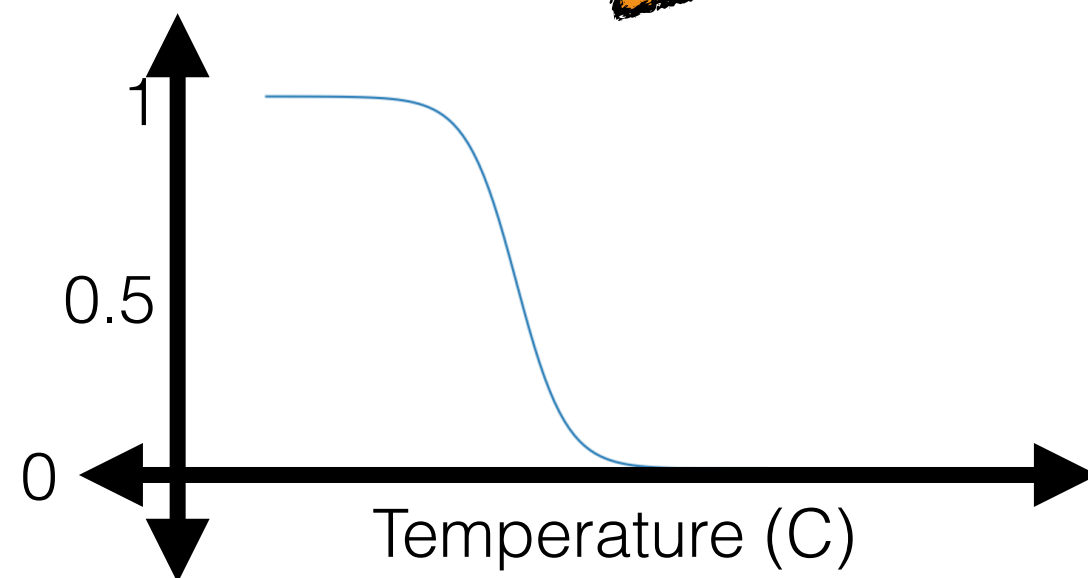- What's an appropriate loss for this guess?

Probability(data)

$$= \prod_{i=1}^{n} \text{Probability}(\text{data point } i)$$

$[\text{Let } g^{(i)} = \sigma(\theta^\top x^{(i)} + \theta_0)]$

$$= \prod_{i=1}^{n} \begin{cases} g^{(i)} \text{ if } y^{(i)} = +1 \\ (1 - g^{(i)}) \text{ else} \end{cases}$$

$$= \prod_{i=1}^{n} (g^{(i)})^{\mathbf{1}\{y^{(i)} = +1\}} (1 - g^{(i)})^{\mathbf{1}\{y^{(i)} \neq +1\}}$$

Loss(data) = -log probability(data)

$$= \frac{1}{n} \sum_{i=1}^{n} - \left( \mathbf{1}\{y^{(i)} = +1\} \log g^{(i)} + \mathbf{1}\{y^{(i)} \neq +1\} \log(1 - g^{(i)}) \right)$$

Temperature (C)

Temperature (C)

+ + +  + + +

− −− − −

9

# Linear logistic classification

- What's an appropriate loss for this guess?

Probability(data)

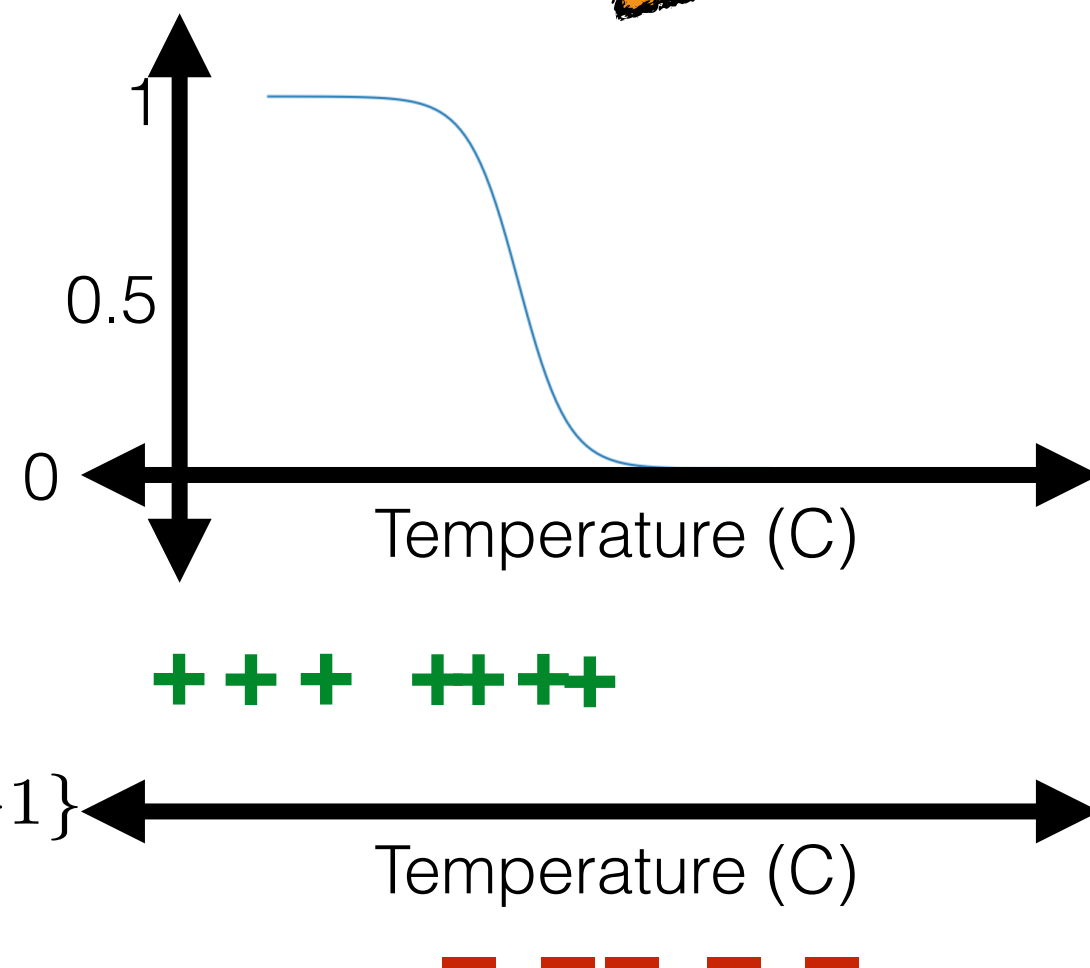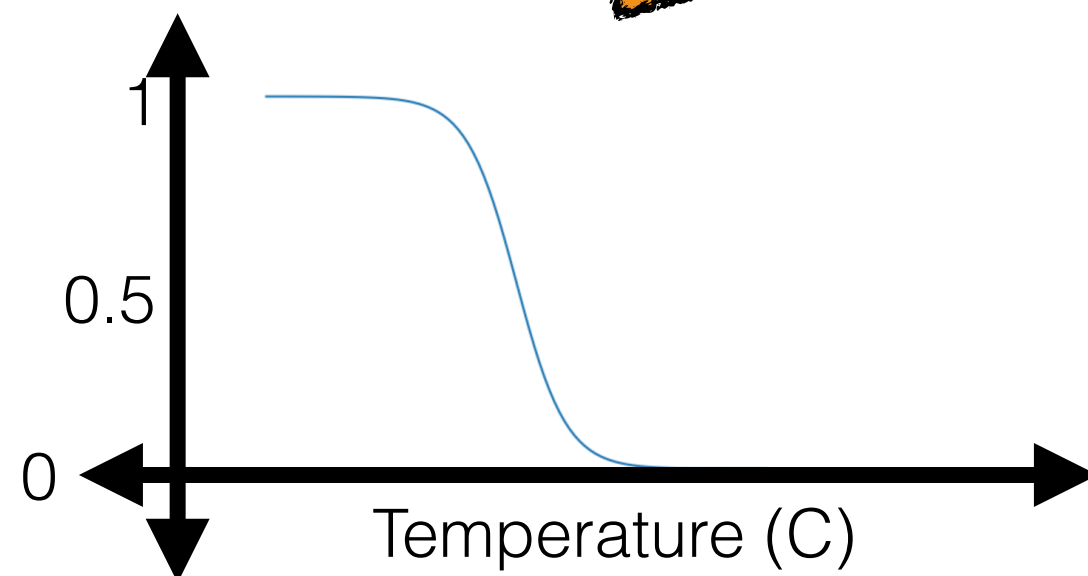$$= \prod_{i=1}^{n} \text{Probability(data point } i)$$

$$[\text{Let } g^{(i)} = \sigma(\theta^{\top} x^{(i)} + \theta_0)]$$

$$= \prod_{i=1}^{n} \begin{cases} g^{(i)} \text{ if } y^{(i)} = +1 \\ (1 - g^{(i)}) \text{ else} \end{cases}$$

$$= \prod_{i=1}^{n} (g^{(i)})^{\mathbf{1}\{y^{(i)}=+1\}} (1 - g^{(i)})^{\mathbf{1}\{y^{(i)} \neq +1\}}$$

Loss(data) = -(1/n) * log probability(data)

$$= \frac{1}{n} \sum_{i=1}^{n} - \left( \mathbf{1}\{y^{(i)} = +1\} \log g^{(i)} + \mathbf{1}\{y^{(i)} \neq +1\} \log(1 - g^{(i)}) \right)$$

1

0.5

0

Temperature (C)

+ + +   + + + +

Temperature (C)

− − − −

9

# Linear logistic classification

- What's an appropriate loss for this guess?

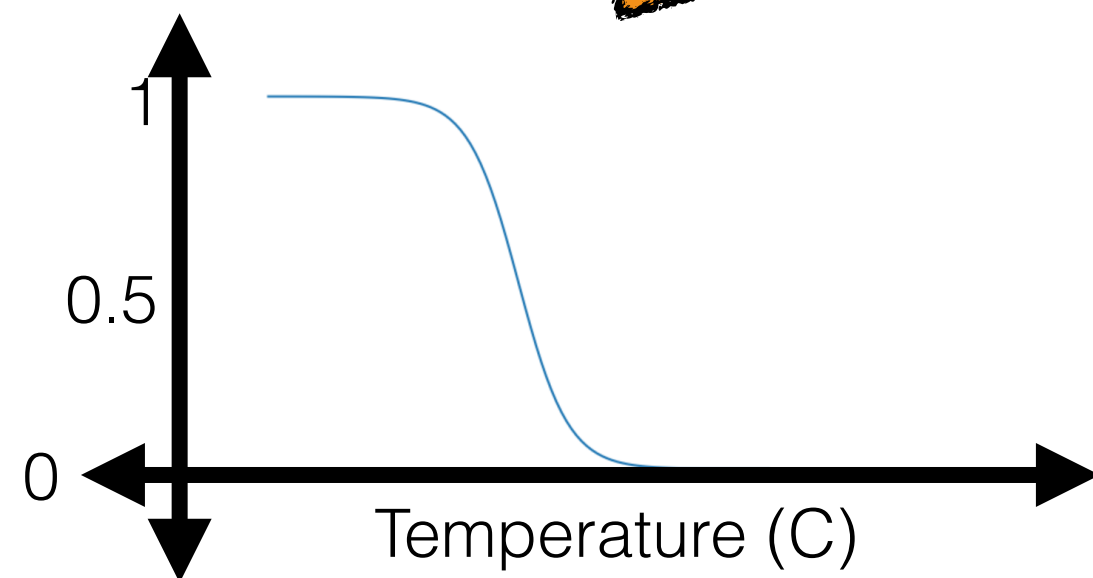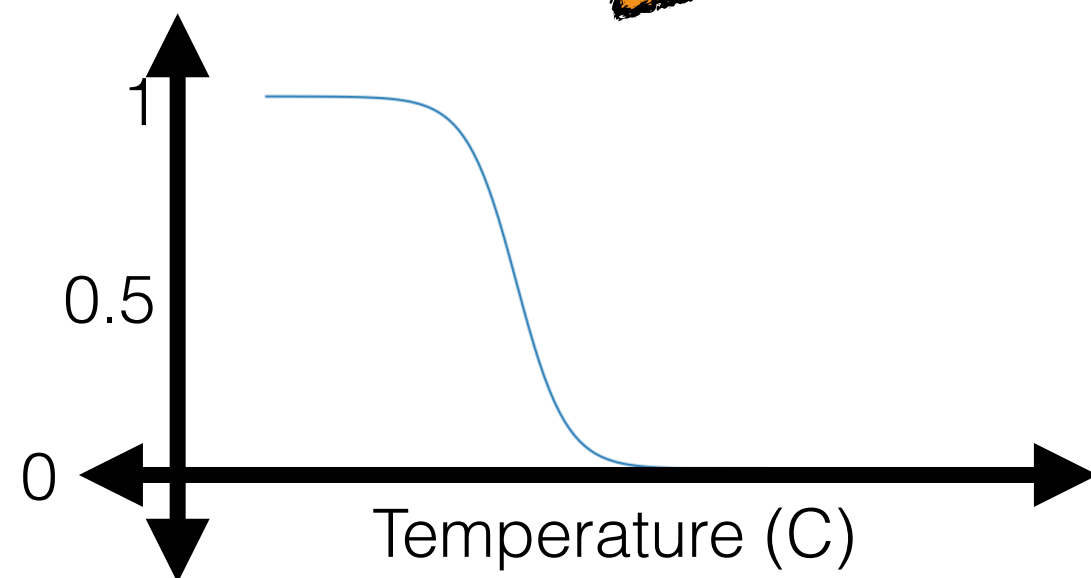Probability(data)

$$= \prod_{i=1}^{n} \text{Probability}(\text{data point } i)$$

$$[\text{Let } g^{(i)} = \sigma(\theta^\top x^{(i)} + \theta_0)]$$

$$= \prod_{i=1}^{n} \begin{cases} g^{(i)} \text{ if } y^{(i)} = +1 \\ (1 - g^{(i)}) \text{ else} \end{cases}$$

$$= \prod_{i=1}^{n} (g^{(i)})^{\mathbf{1}\{y^{(i)}=+1\}} (1 - g^{(i)})^{\mathbf{1}\{y^{(i)}\neq+1\}}$$

Loss(data) = -(1/n) * log probability(data)

$$= \frac{1}{n} \sum_{i=1}^{n} -\left( \mathbf{1}\{y^{(i)} = +1\} \log g^{(i)} + \mathbf{1}\{y^{(i)} \neq +1\} \log(1 - g^{(i)}) \right)$$

Negative log likelihood loss (*g* for guess, *a* for actual):

$$-L_{\text{nll}}(g, a) = (1\{a = +1\} \log g + 1\{a \neq +1\} \log(1 - g))$$

9

# Gradient descent for logistic regression

- Want to minimize average (negative log likelihood) loss across the data (objective is differentiable and convex)
  $$J_{\mathrm{lr}}(\Theta) = J_{\mathrm{lr}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} L_{\mathrm{nll}}(\sigma(\theta^\top x^{(i)} + \theta_0), y^{(i)})$$

- Run `Gradient-Descent(`$\Theta_{\mathrm{init}}, \eta, J_{lr}, \nabla_\Theta J_{lr}, \epsilon$`)`

# Gradient descent for logistic regression

- Want to minimize average (negative log likelihood) loss across the data (objective is differentiable and convex)

$$J_{\mathrm{lr}}(\Theta) = J_{\mathrm{lr}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} L_{\mathrm{nll}}(\sigma(\theta^\top x^{(i)} + \theta_0), y^{(i)})$$

- Run `Gradient-Descent(`$\Theta_{\mathrm{init}}, \eta, J_{lr}, \nabla_\Theta J_{lr}, \epsilon$`)`
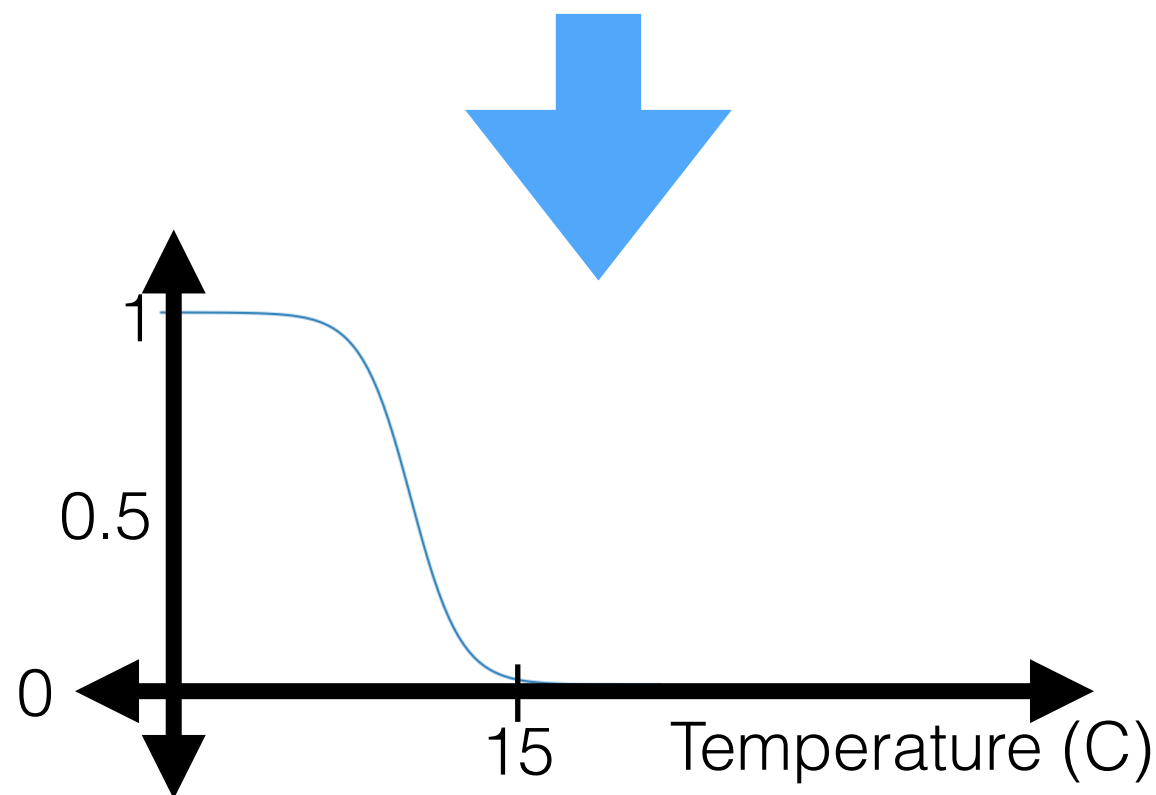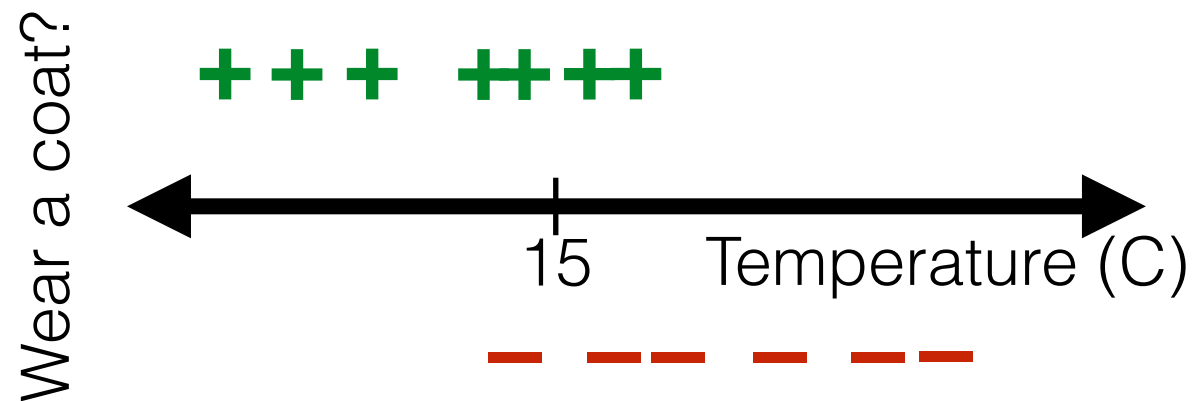
# Gradient descent for logistic regression

- Want to minimize average (negative log likelihood) loss across the data (objective is differentiable and convex)

$$J_{\mathrm{lr}}(\Theta) = J_{\mathrm{lr}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} L_{\mathrm{nll}}(\sigma(\theta^\top x^{(i)} + \theta_0), y^{(i)})$$

- Run Gradient-Descent($\Theta_{\mathrm{init}}, \eta, J_{lr}, \nabla_\Theta J_{lr}, \epsilon$)
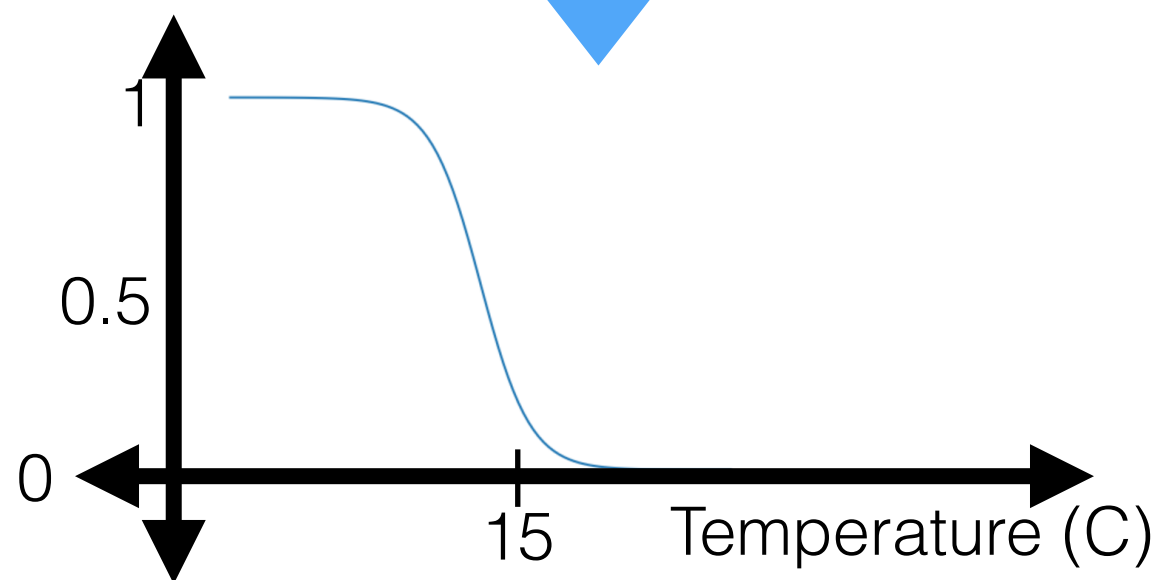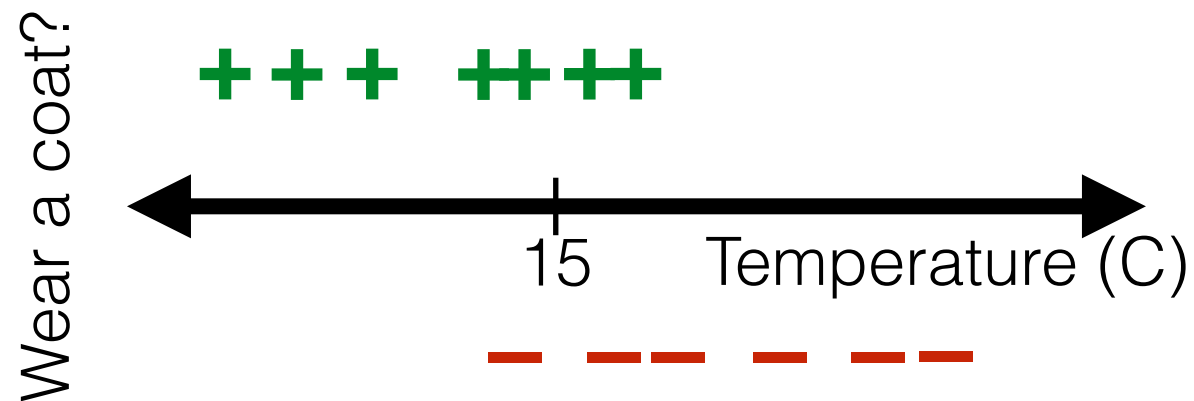
# Gradient descent for logistic regression

- Want to minimize average (negative log likelihood) loss across the data (objective is differentiable and convex)

$$J_{\mathrm{lr}}(\Theta) = J_{\mathrm{lr}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} L_{\mathrm{nll}}(\sigma(\theta^\top x^{(i)} + \theta_0), y^{(i)})$$

- Run `Gradient-Descent(`$\Theta_{\mathrm{init}}, \eta, J_{lr}, \nabla_\Theta J_{lr}, \epsilon$`)`
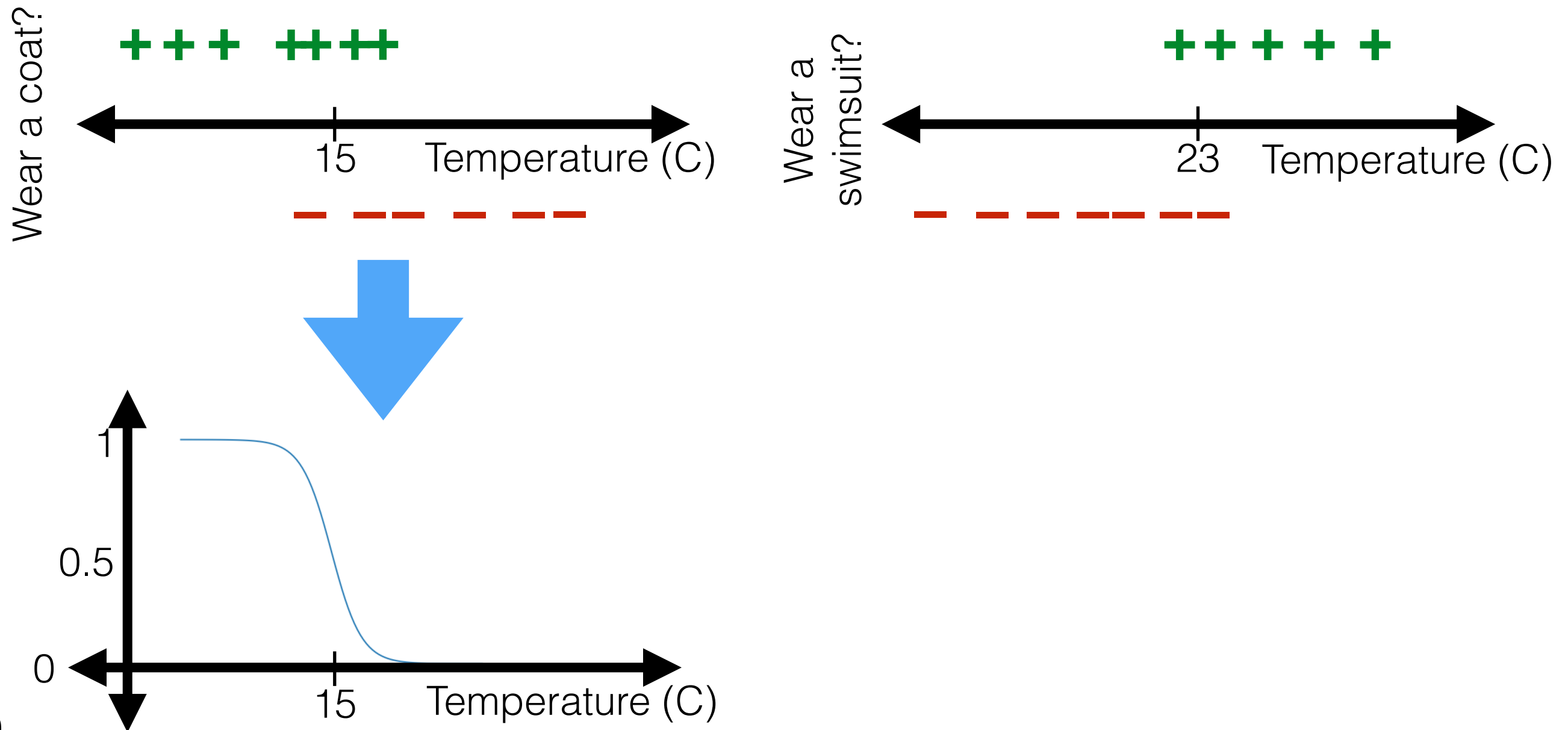
# Gradient descent for logistic regression

- Want to minimize average (negative log likelihood) loss across the data (objective is differentiable and convex)

$$J_{\mathrm{lr}}(\Theta) = J_{\mathrm{lr}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} L_{\mathrm{nll}}(\sigma(\theta^\top x^{(i)} + \theta_0), y^{(i)})$$

- Run Gradient-Descent( $\Theta_{\mathrm{init}}, \eta, J_{lr}, \nabla_\Theta J_{lr}, \epsilon$ )
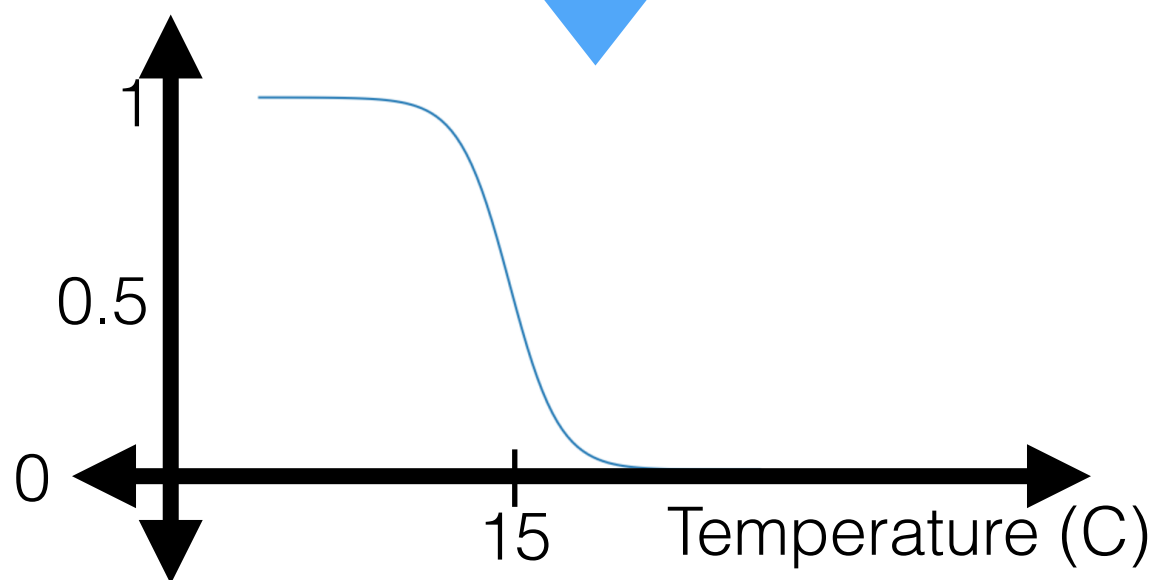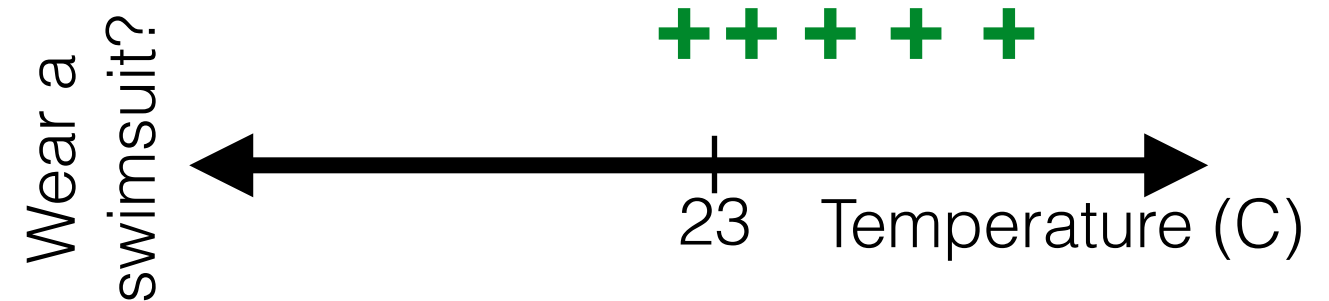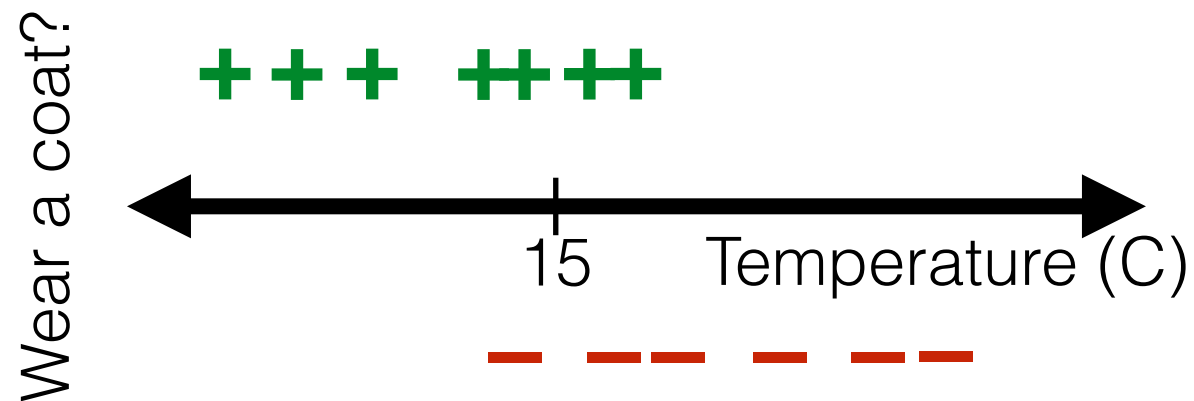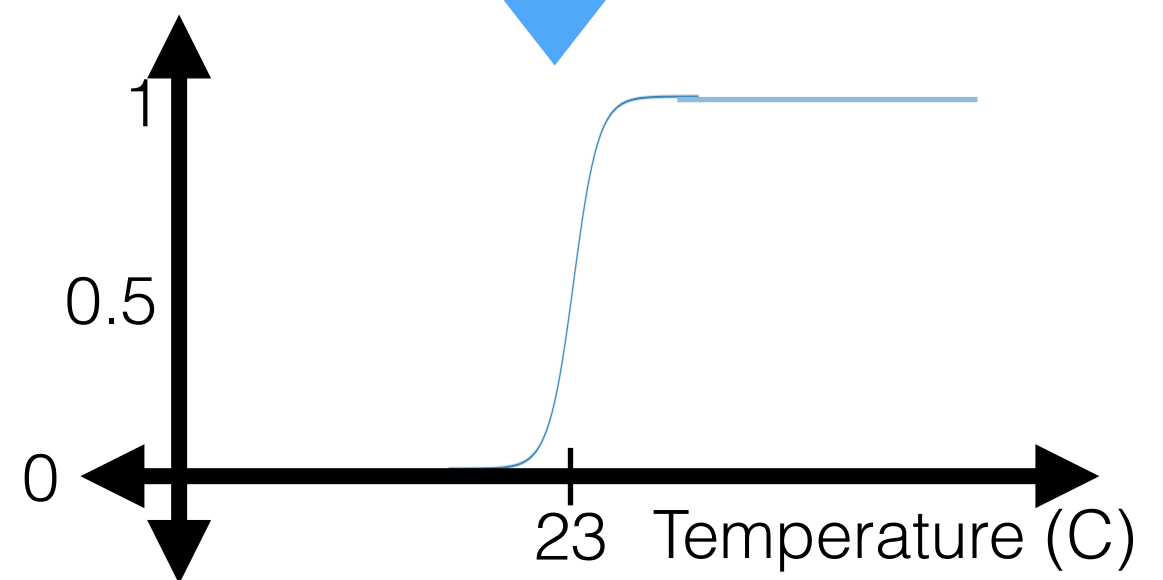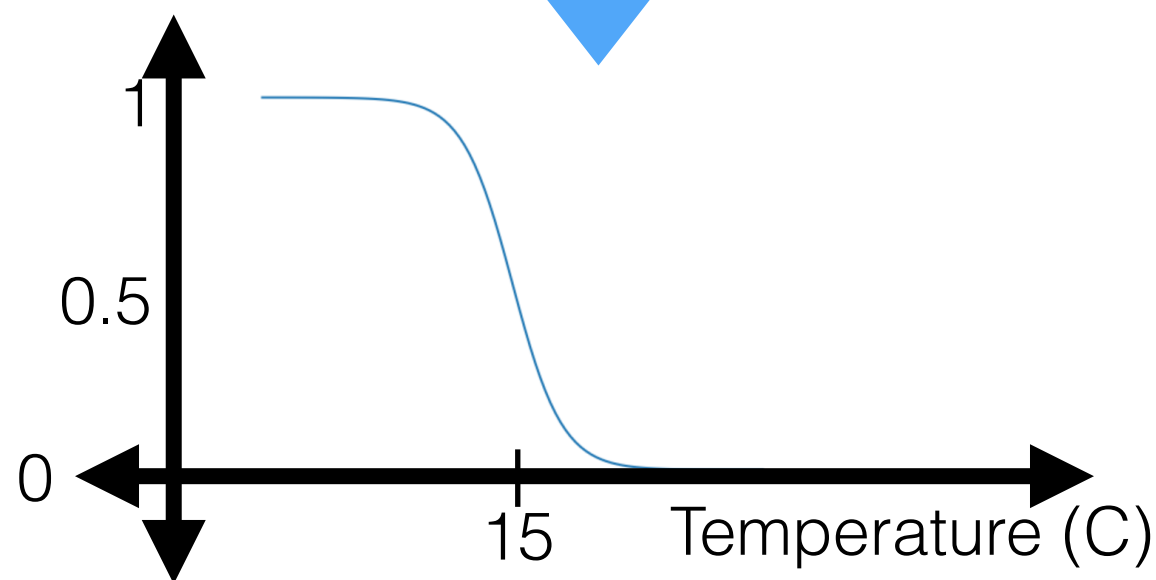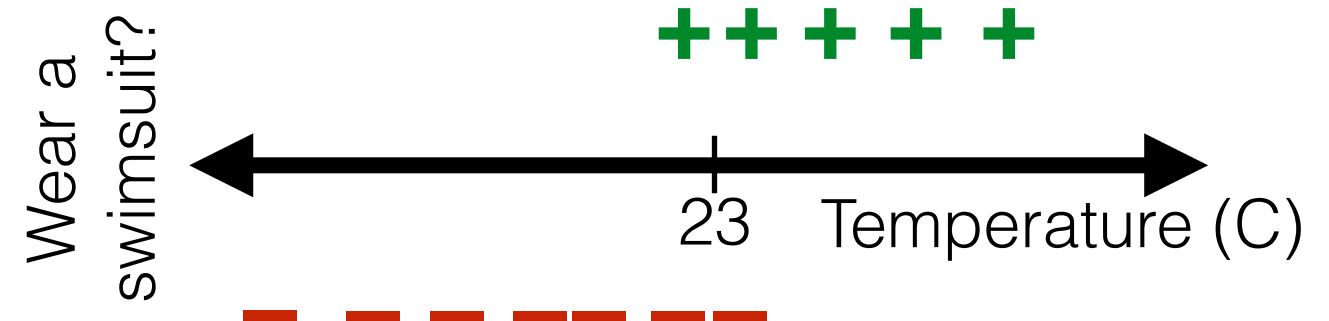
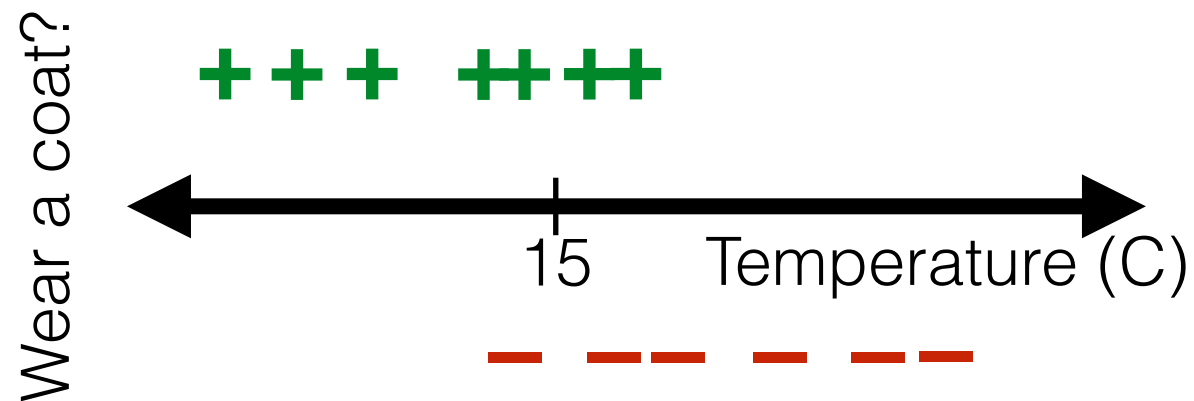# Gradient descent for logistic regression

- Want to minimize average (negative log likelihood) loss across the data (objective is differentiable and convex)
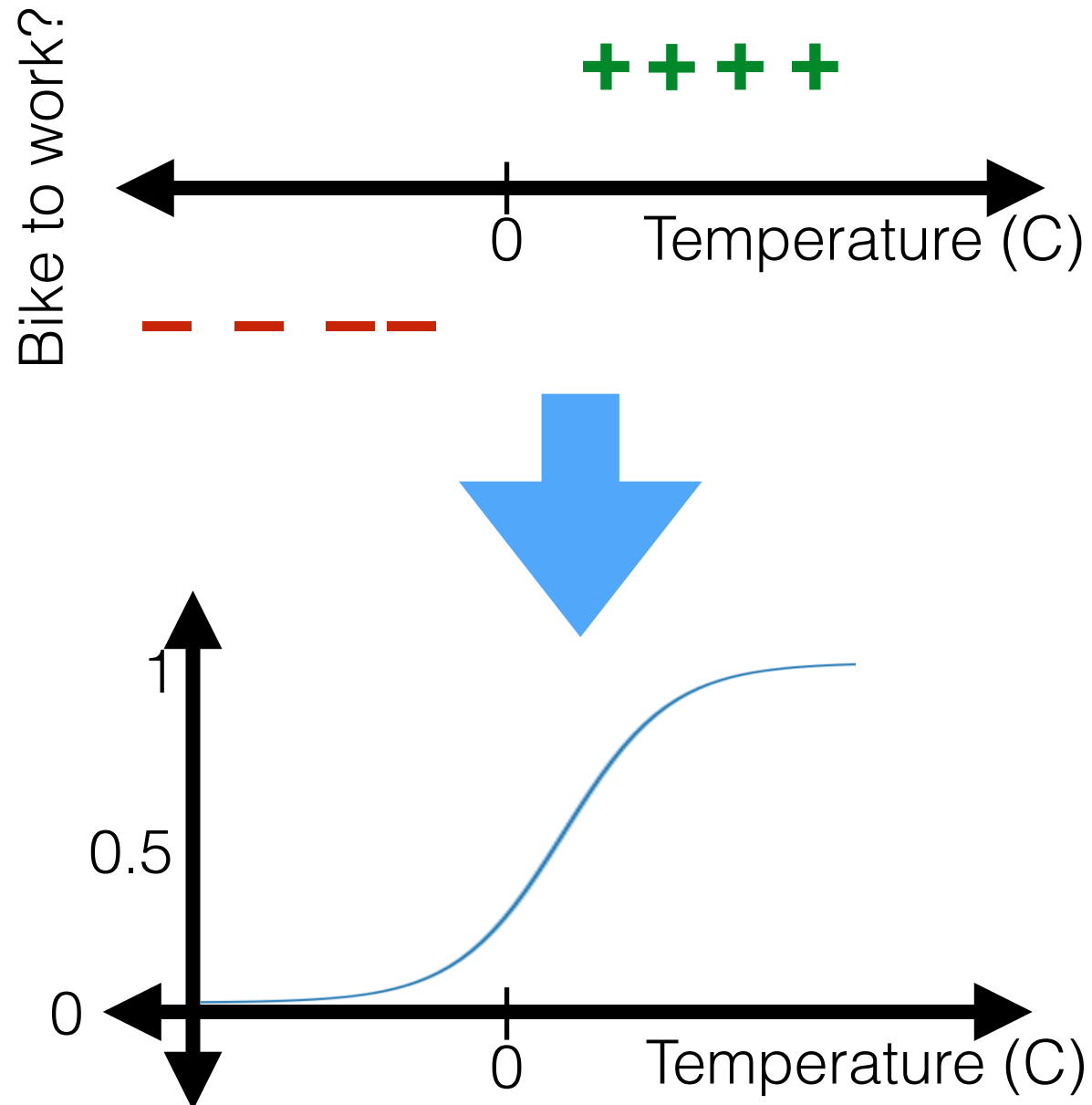$$J_{\mathrm{lr}}(\Theta) = J_{\mathrm{lr}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} L_{\mathrm{nll}}(\sigma(\theta^\top x^{(i)} + \theta_0), y^{(i)})$$

- Run `Gradient-Descent(`$\Theta_{\mathrm{init}}, \eta, J_{lr}, \nabla_\Theta J_{lr}, \epsilon$`)`

# Gradient descent for logistic regression

- Can still have practical issues though!

- Run `Gradient-Descent(`$\Theta_{\text{init}}, \eta, J_{lr}, \nabla_{\Theta}J_{lr}, \epsilon$`)`

# Gradient descent for logistic regression

- Can still have practical issues though!

- Run `Gradient-Descent(` $\Theta_{\text{init}}, \eta, J_{lr}, \nabla_\Theta J_{lr}, \epsilon$ `)`
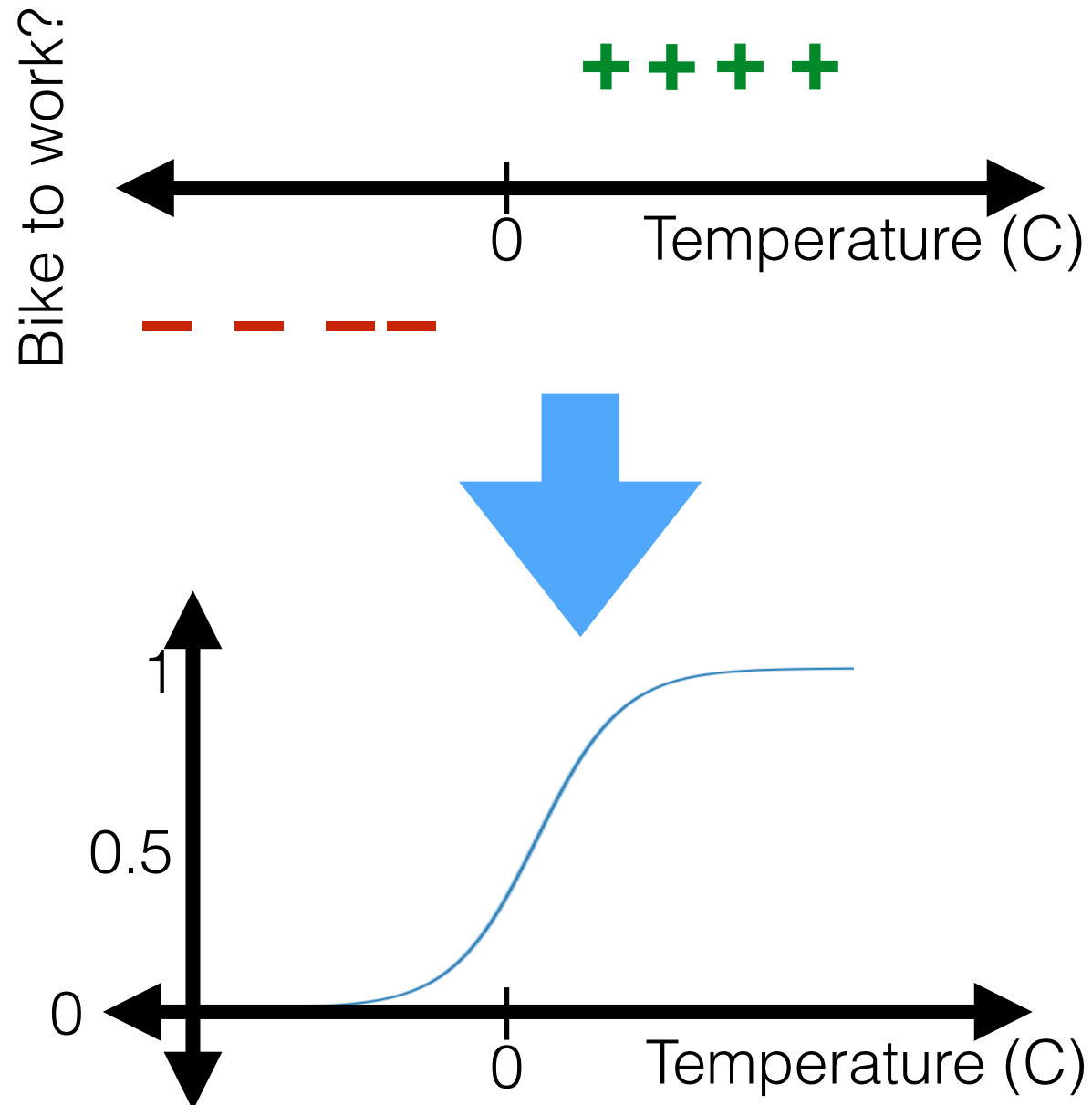
# Gradient descent for logistic regression

- Can still have practical issues though!

- Run `Gradient-Descent(`$\Theta_{\text{init}}, \eta, J_{lr}, \nabla_\Theta J_{lr}, \epsilon$`)`
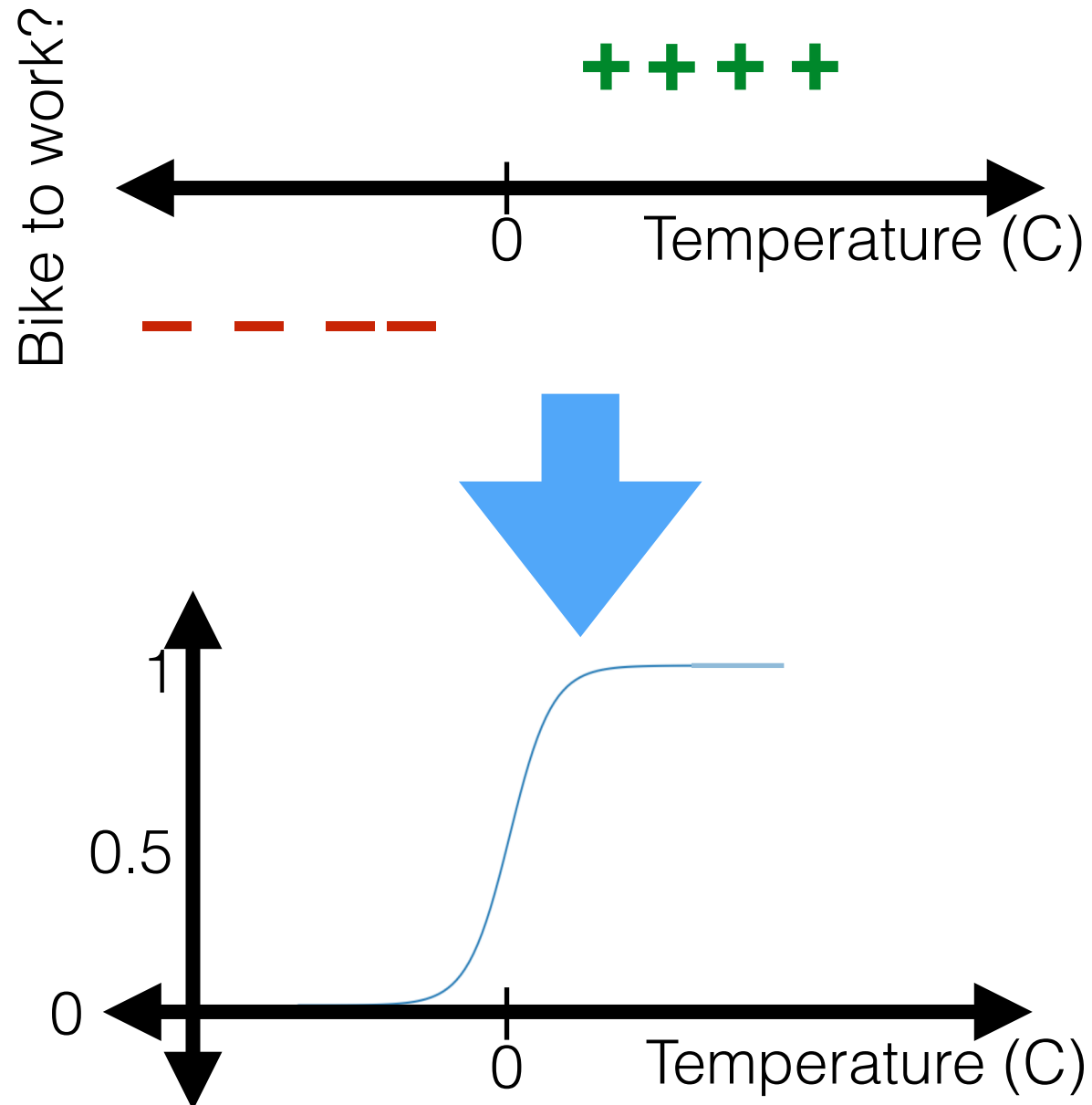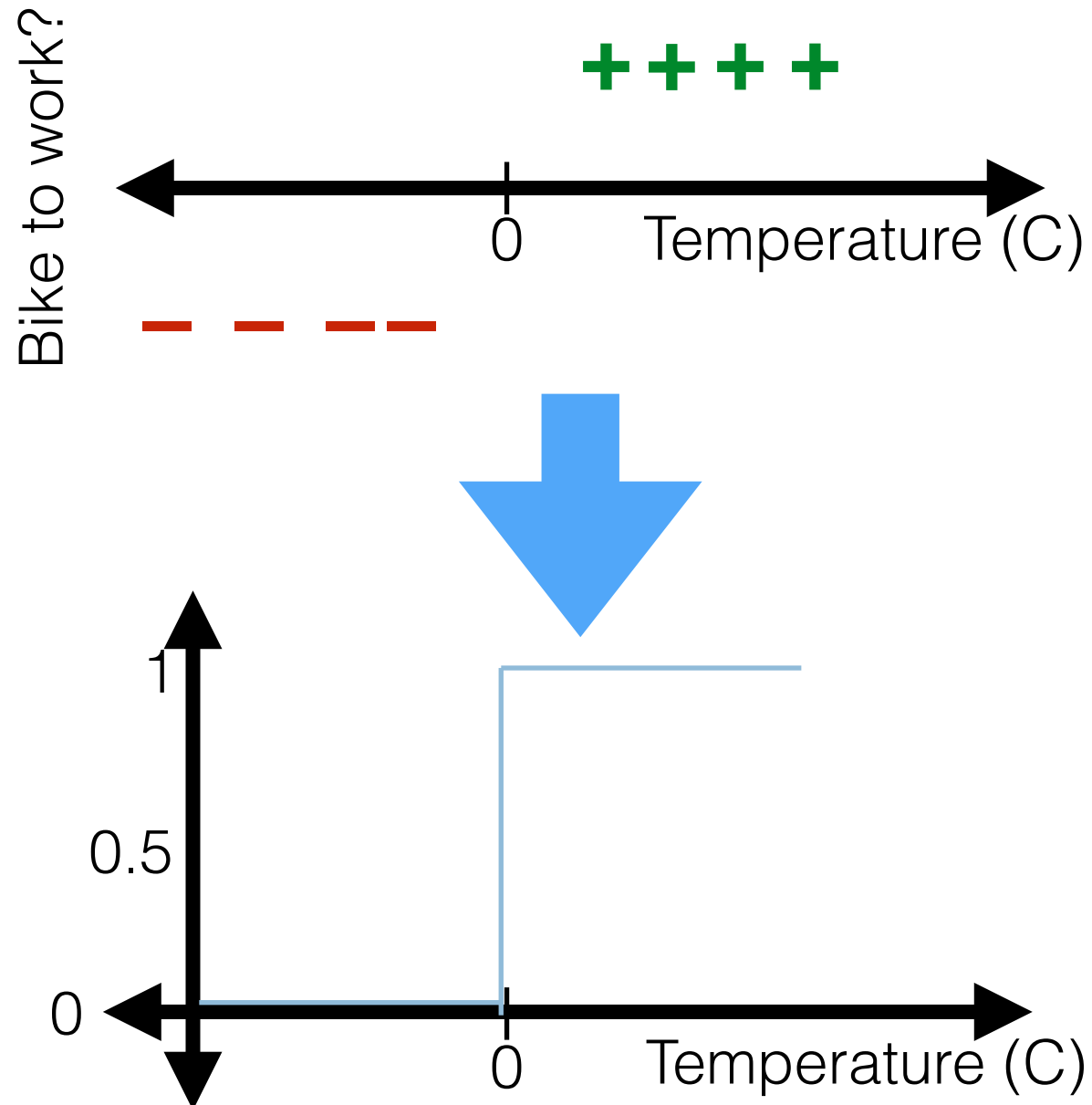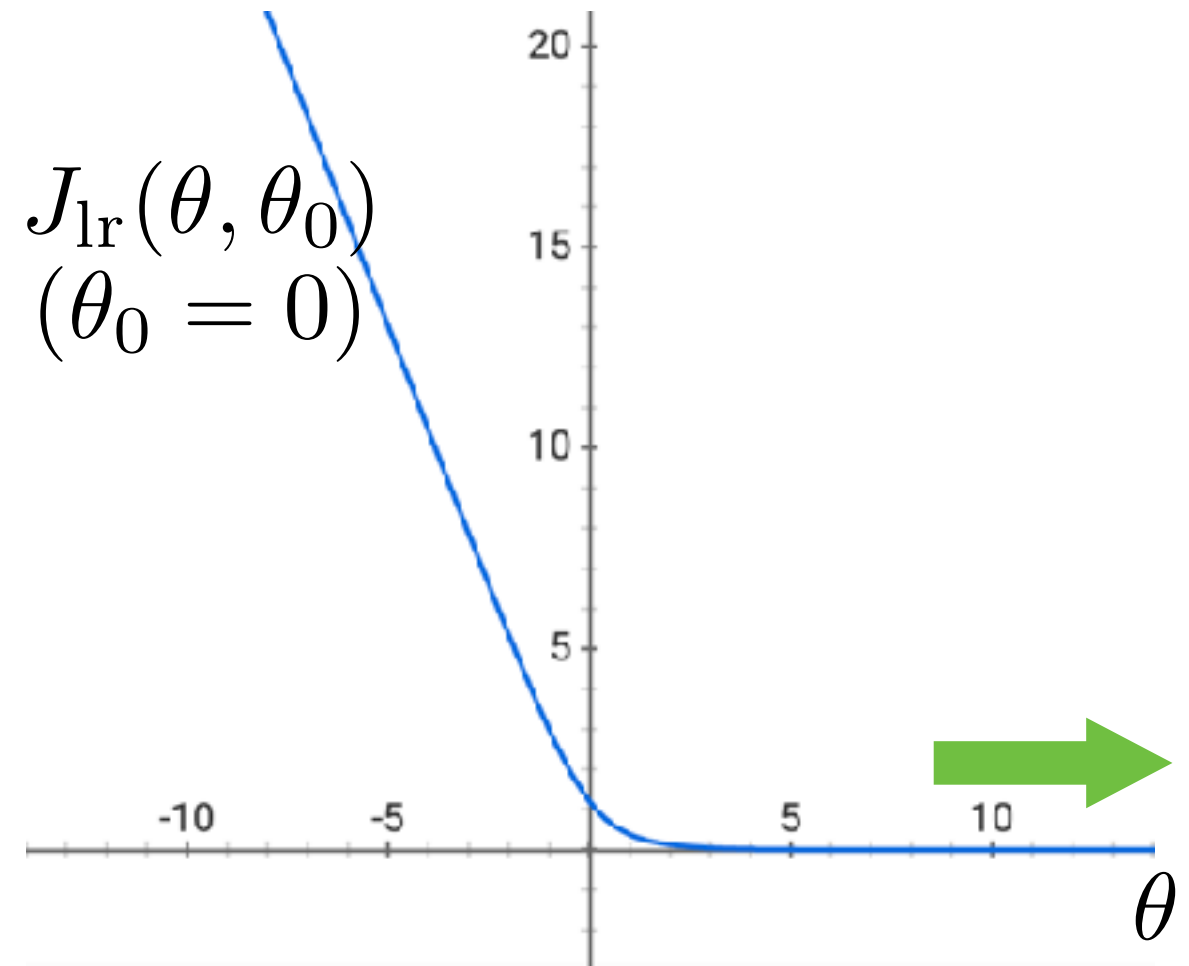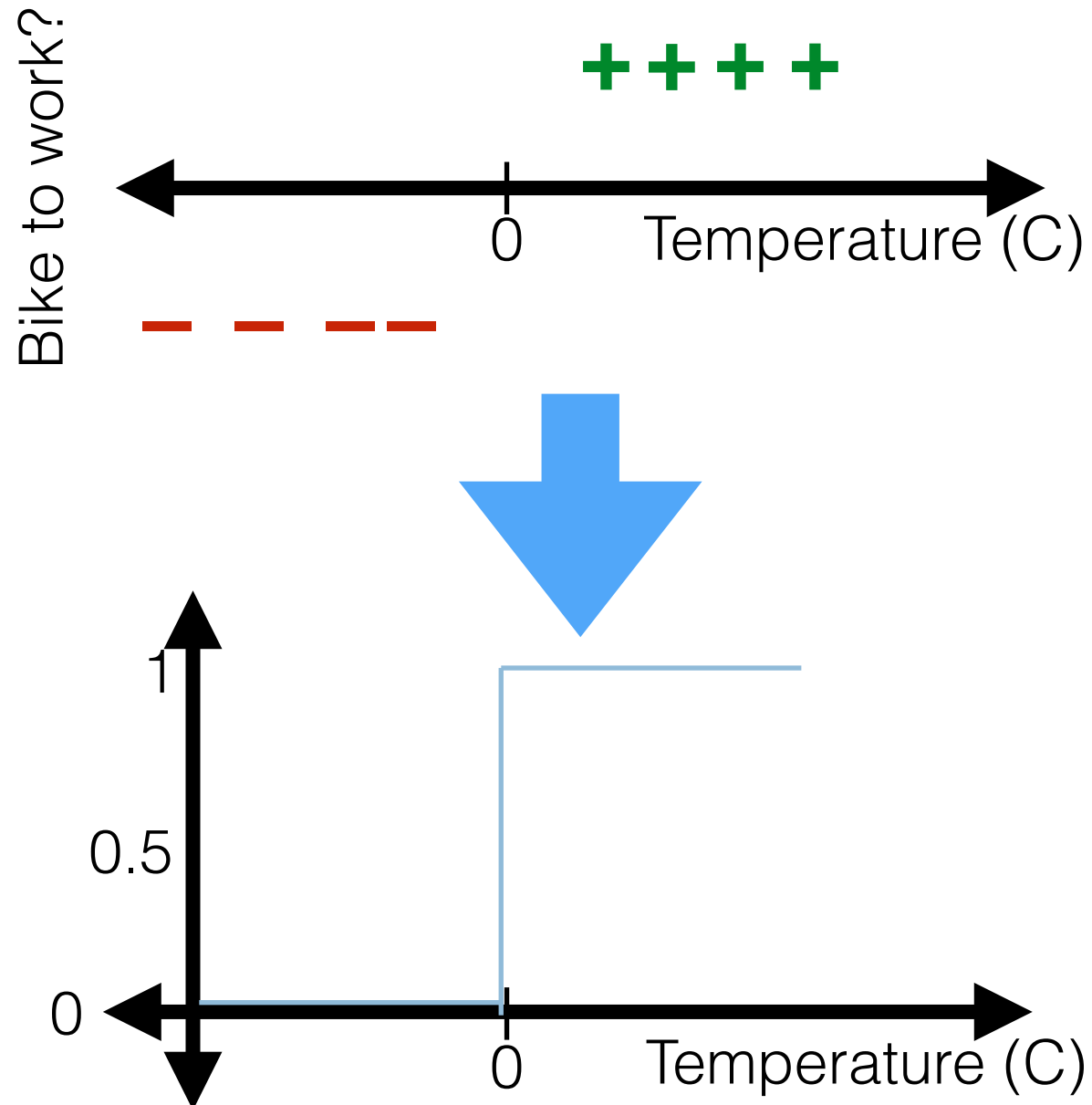
# Gradient descent for logistic regression

- Can still have practical issues though!
- Run `Gradient-Descent(` $\Theta_{\mathrm{init}}, \eta, J_{lr}, \nabla_\Theta J_{lr}, \epsilon$ `)`

# Gradient descent for logistic regression

- Can still have practical issues though!
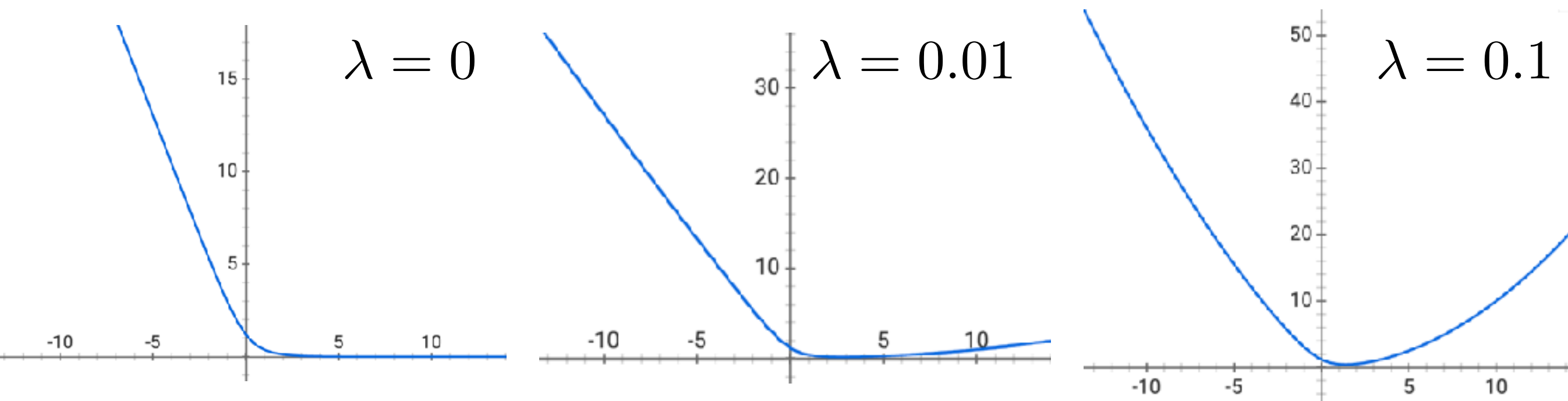- Run $\texttt{Gradient-Descent}(\Theta_{\text{init}}, \eta, J_{lr}, \nabla_{\Theta} J_{lr}, \epsilon)$



$J_{\text{lr}}(\theta, \theta_0)$
$(\theta_0 = 0)$

no global optimum

# Logistic regression loss revisited

$$J_{\mathrm{lr}}(\Theta) = J_{\mathrm{lr}}(\theta, \theta_0)$$

$$= \frac{1}{n} \sum_{i=1}^{n} L_{\mathrm{nll}}(\sigma(\theta^\top x^{(i)} + \theta_0), y^{(i)}) + \lambda \|\theta\|^2 \quad (\lambda \geq 0)$$

- A "regularizer" or "penalty"  $R(\theta) = \lambda \|\theta\|^2$
- Penalizes being overly certain
- Objective is still differentiable & convex (gradient descent)



$\lambda = 0$     $\lambda = 0.01$     $\lambda = 0.1$

- How to choose hyperparameter? One option: consider a handful of possible values and compare via CV