

An abstract graphic on the left side of the slide, featuring a sphere-like shape composed of a dense grid of intersecting red, green, and blue lines. The lines are curved and follow the contour of the sphere, creating a complex, woven pattern. The sphere is set against a dark gray background.

10-315

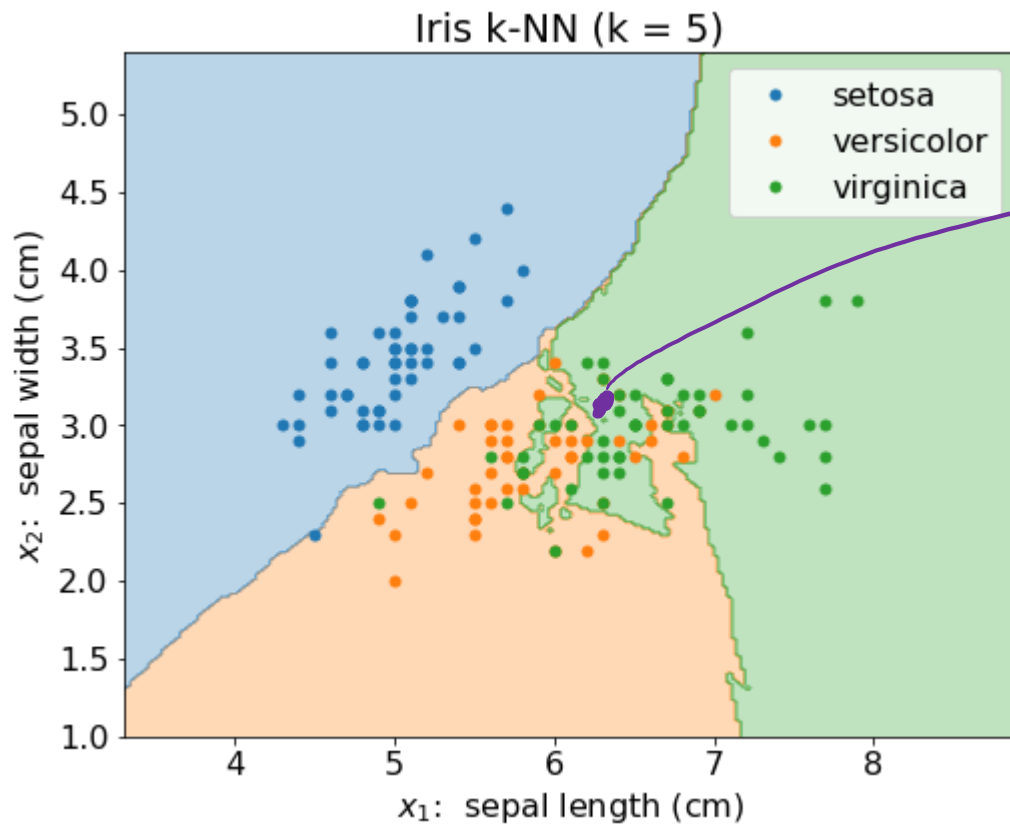
Introduction to ML

Logistic Regression

Instructor: Pat Virtue

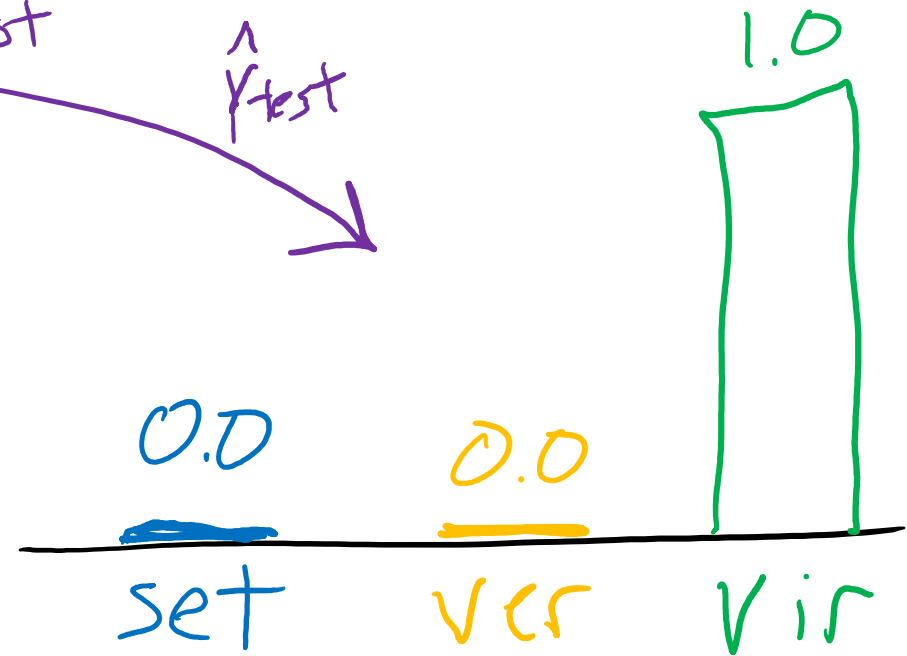
Classification Decisions

Predicting one specific class is troubling, especially when we know that there is some uncertainty in our prediction



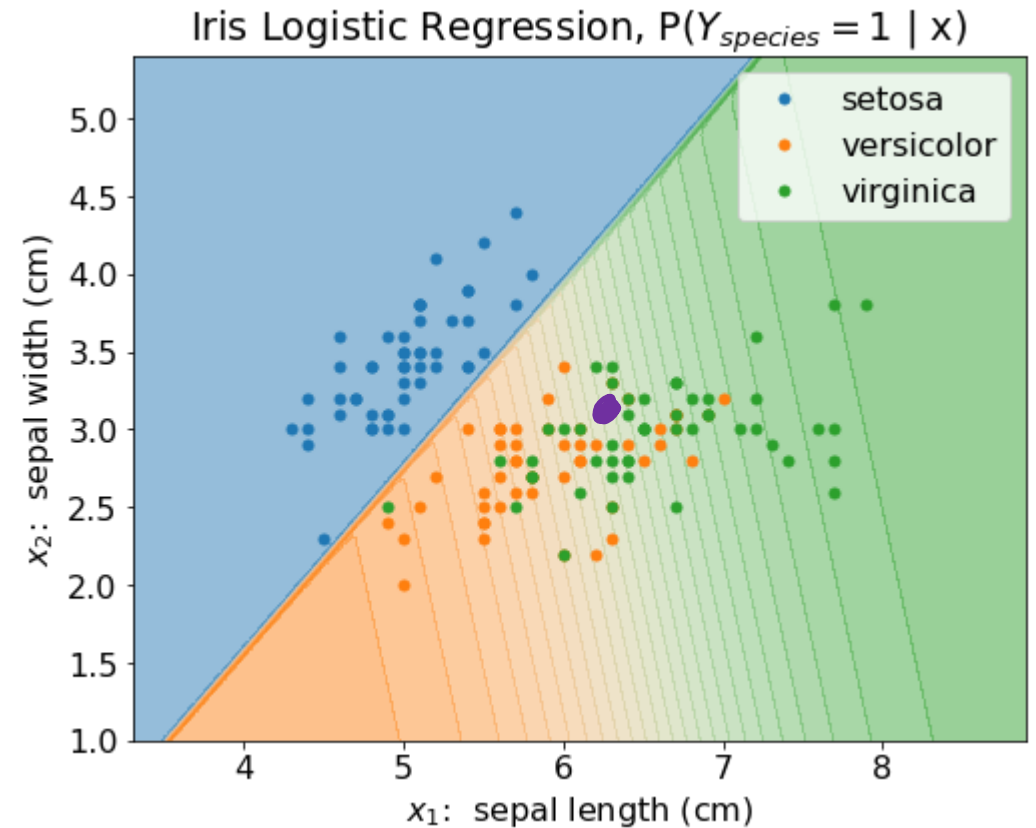
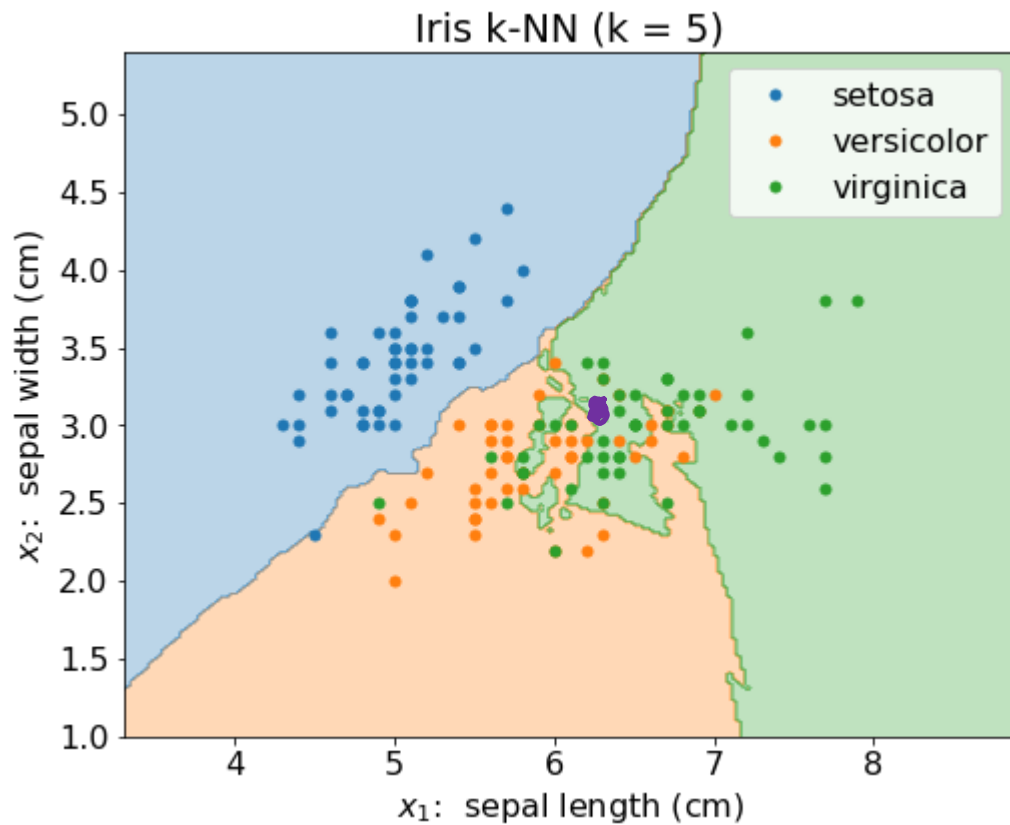
\vec{x}_{test}

\hat{y}_{test}



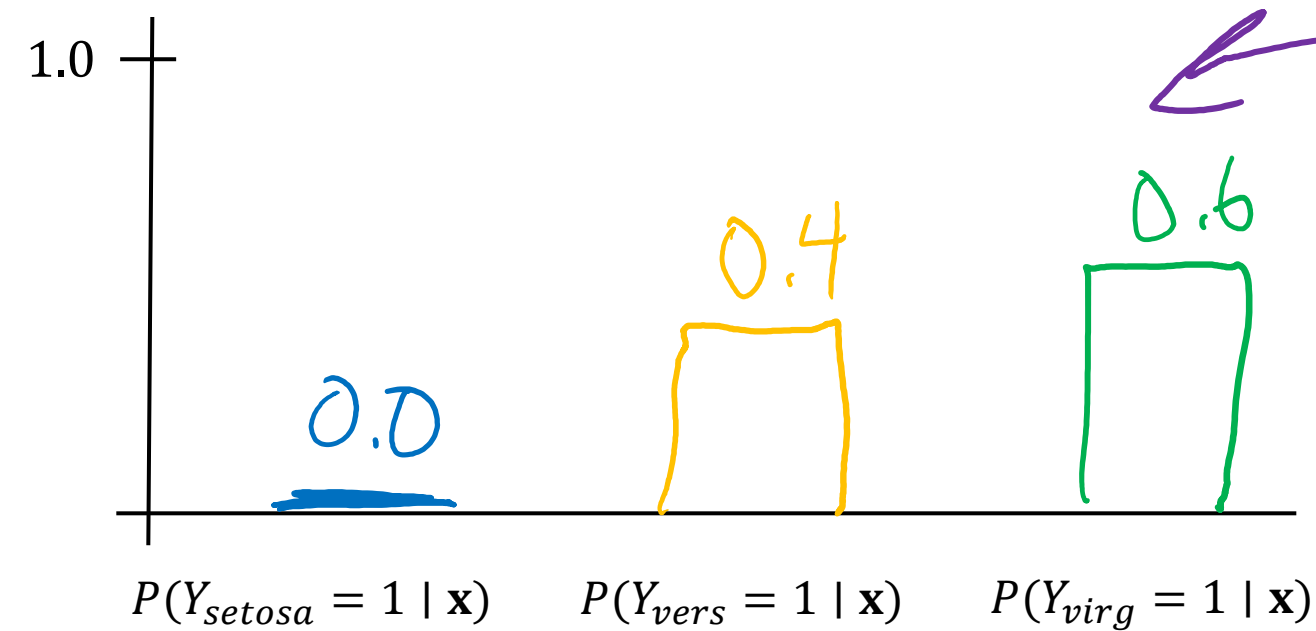
Classification Probability

Constructing a model than can return the probability of the output being a specific class could be incredibly useful

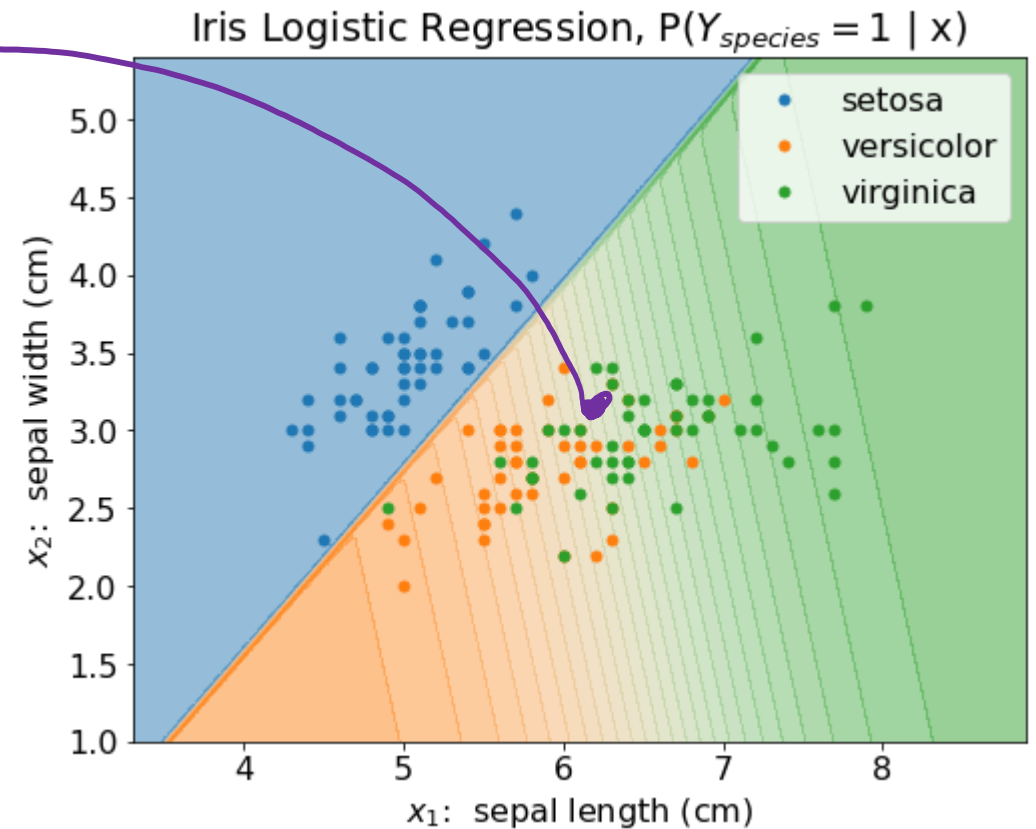


Classification Probability

Constructing a model that can return the probability of the output being a specific class could be incredibly useful

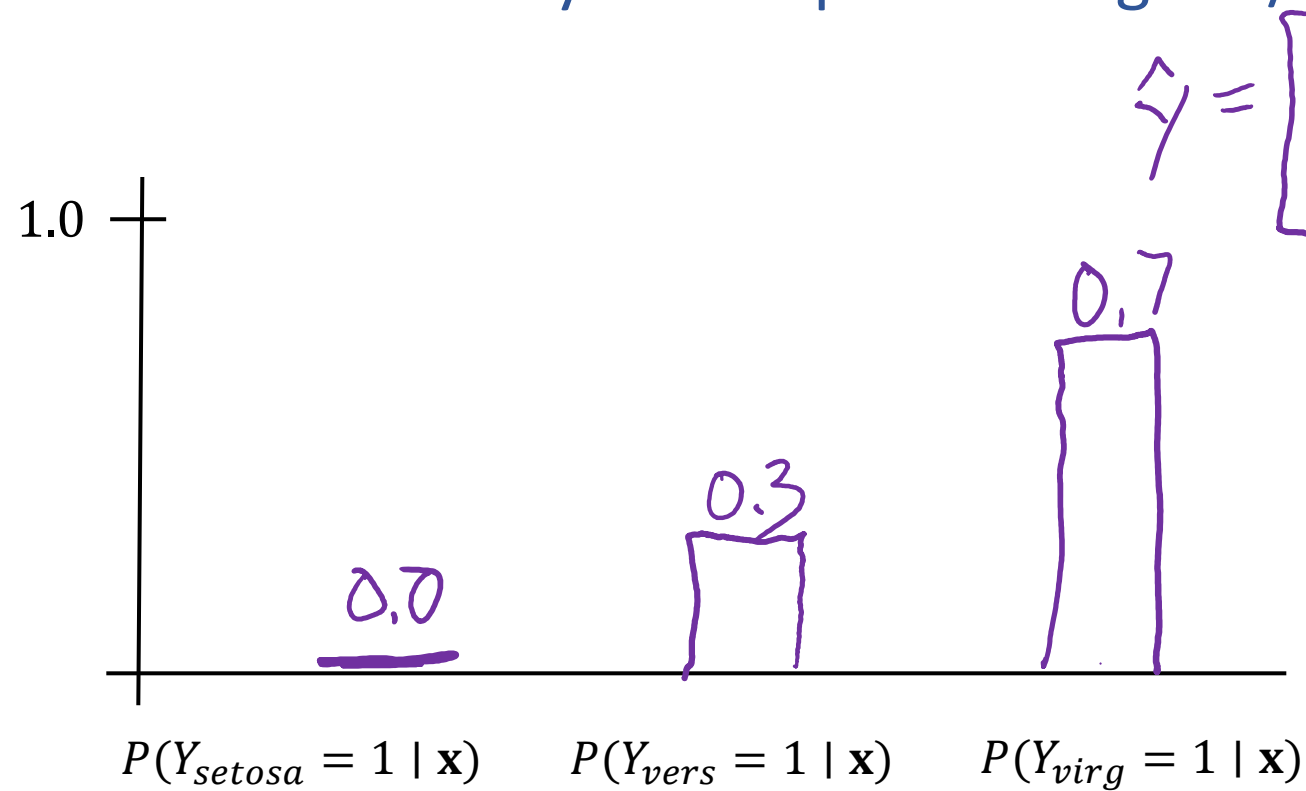


We can still make decisions, .e.g,
$$\operatorname{argmax}_k P(Y_k = 1 | \mathbf{x})$$

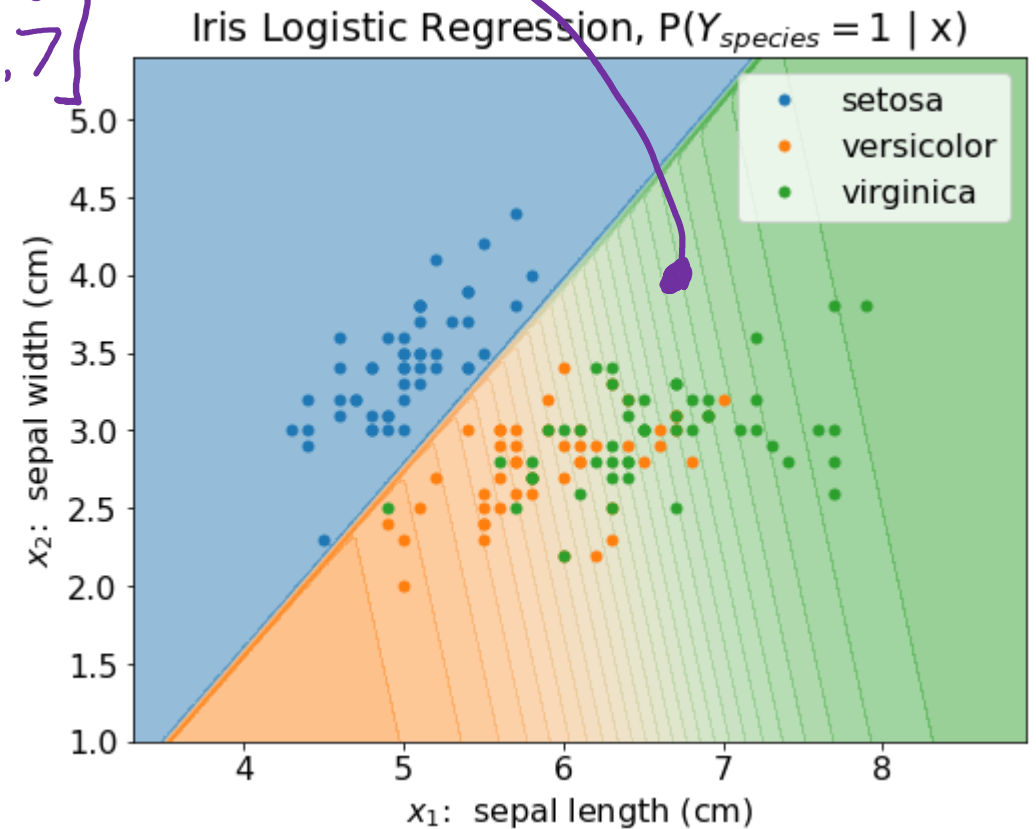


Loss for Probability Distributions

We need a way to compare how good/bad each prediction is



$$\hat{\mathbf{y}} = \begin{bmatrix} 0.0 \\ 0.3 \\ 0.7 \end{bmatrix}$$

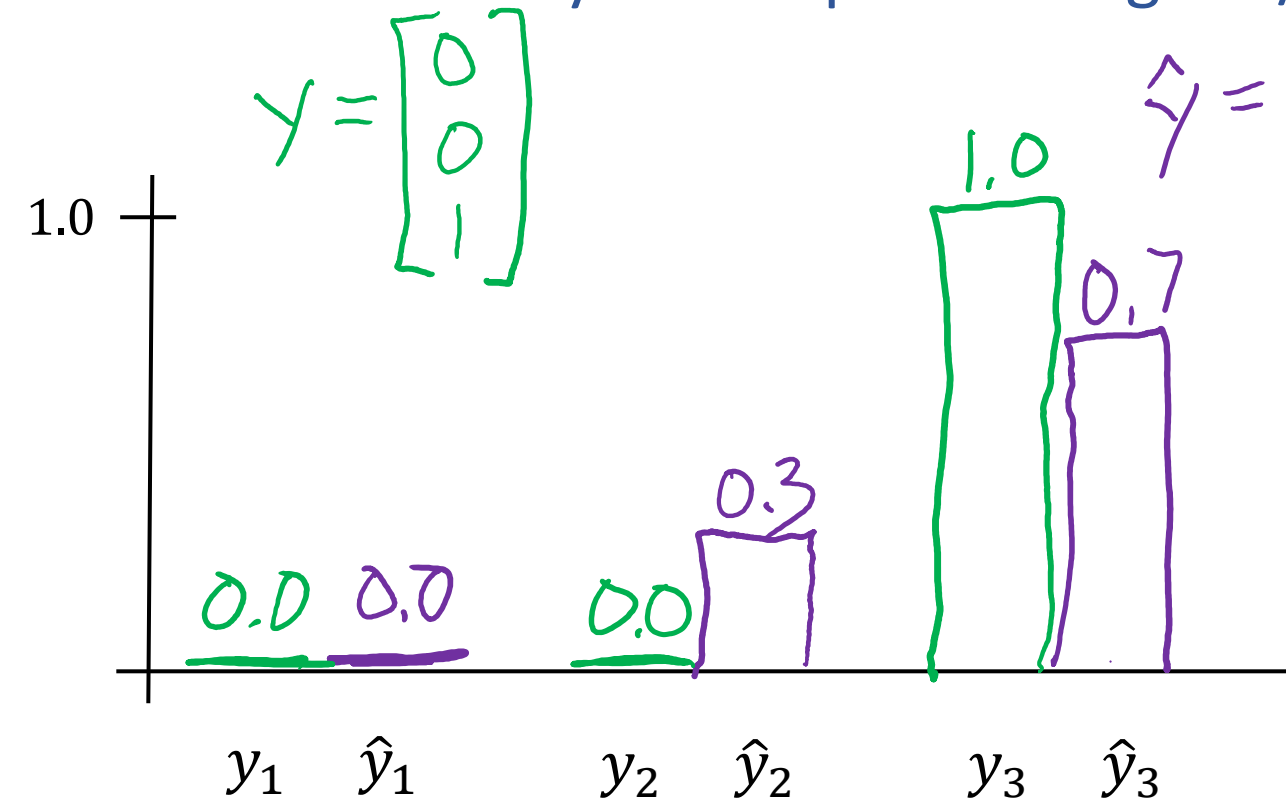


Cross-entropy loss

$$\ell(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{k=1}^K y_k \log \hat{y}_k$$

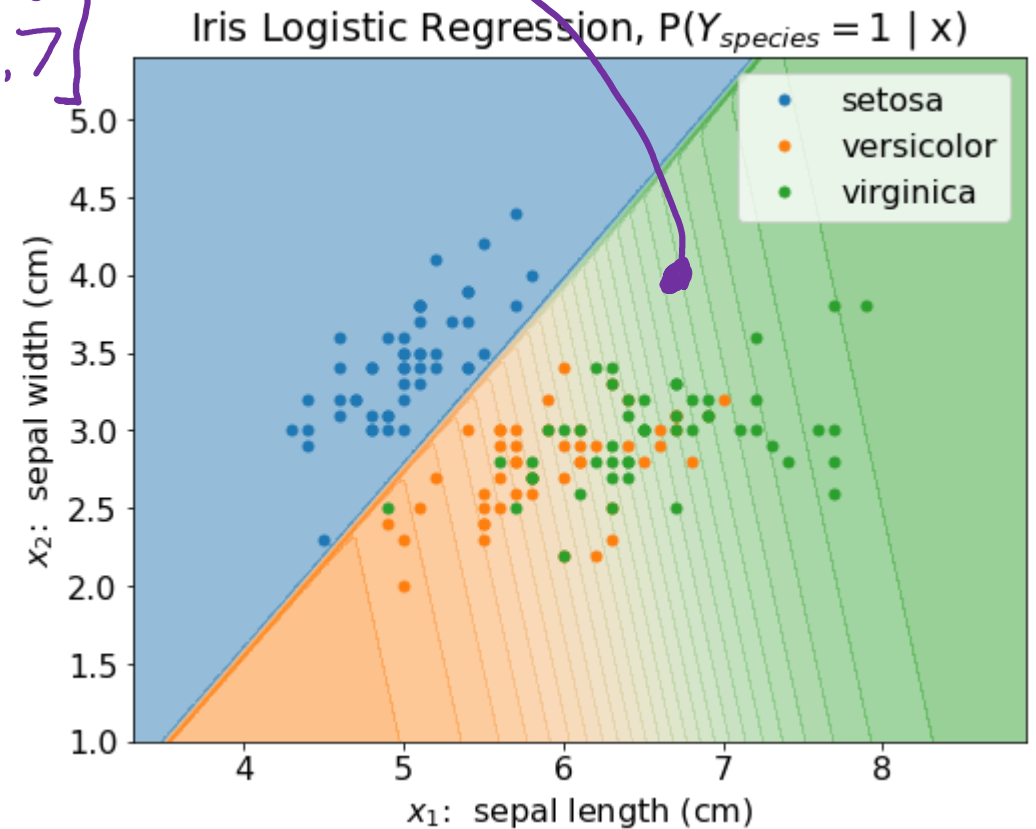
Loss for Probability Distributions

We need a way to compare how good/bad each prediction is



Cross-entropy loss

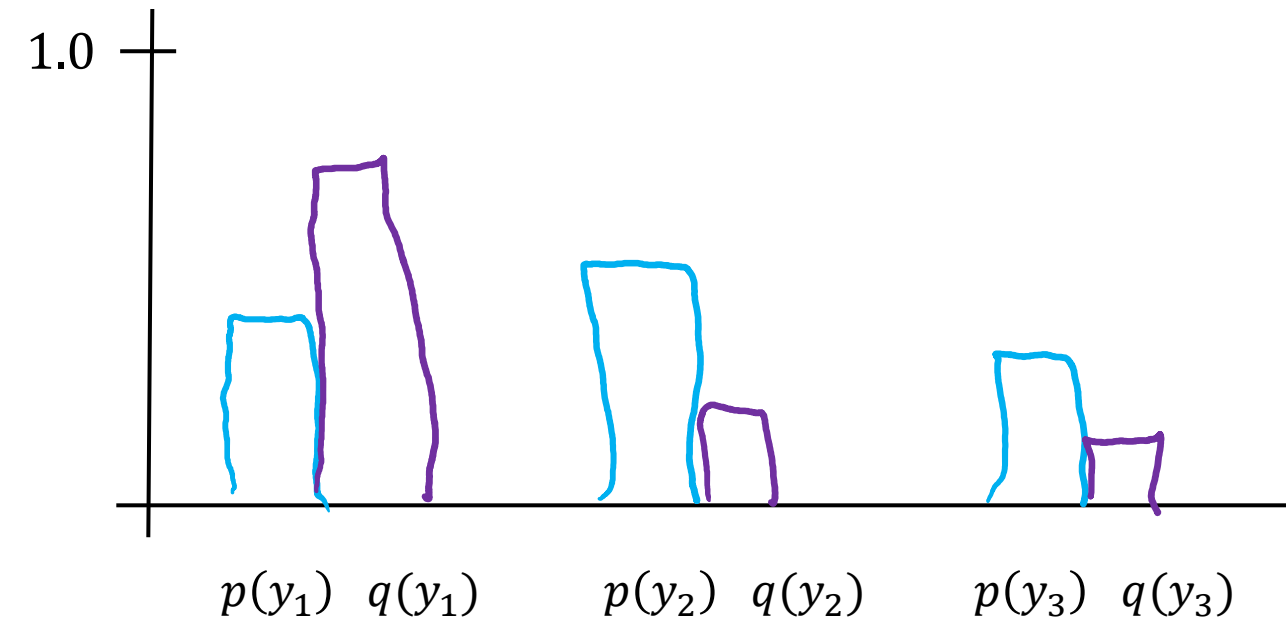
$$\ell(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{k=1}^K \mathbf{y}_k \log \hat{\mathbf{y}}_k$$



Loss for Probability Distributions

Cross-entropy more generally is a way to compare any two probability distributions*

*when used in logistic regression \mathbf{y} is always a one-hot vector



Cross-entropy loss

$$H(\underline{P}, \underline{Q}) = - \sum_{k=1}^K \underline{p(y_k)} \log \underline{q(y_k)}$$

Empirical Risk Minimization

Still doing empirical risk minimization, just with a cross-entropy loss

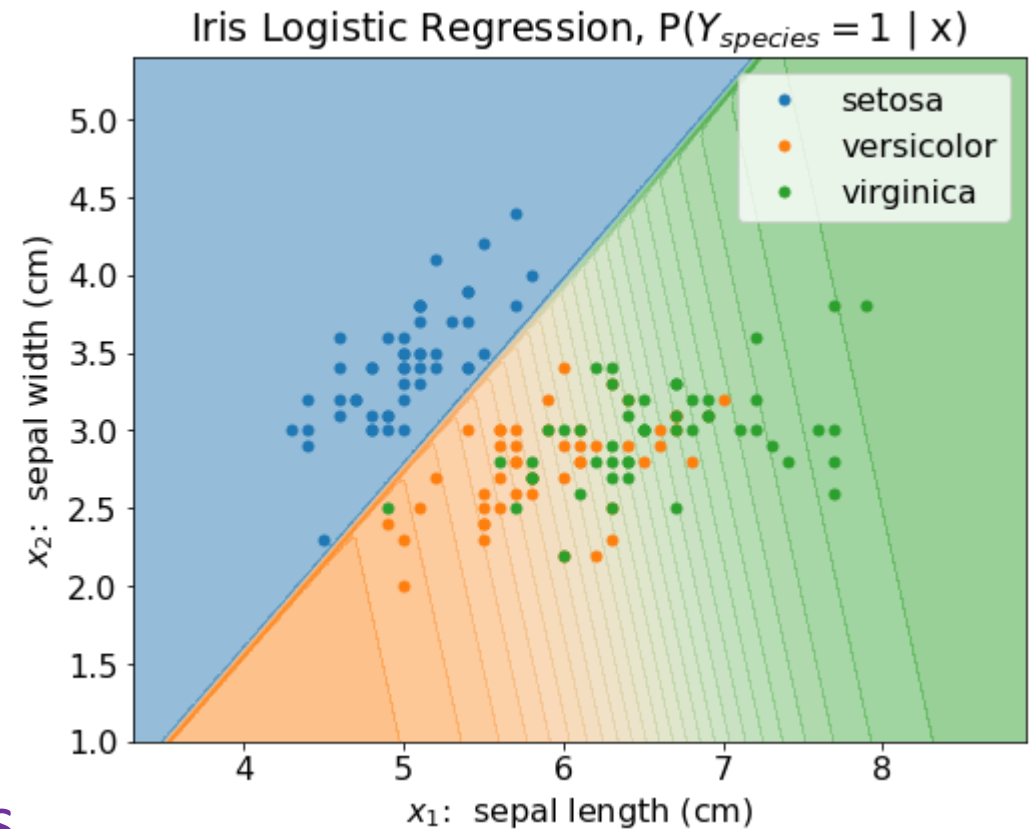
$$h^* = \operatorname{argmin}_{h \in \mathcal{H}} \hat{R}(h)$$

$$\hat{R}(h) = \frac{1}{N} \sum_{i=1}^N \ell \left(y^{(i)}, h \left(x^{(i)} \right) \right)$$

Cross-entropy loss

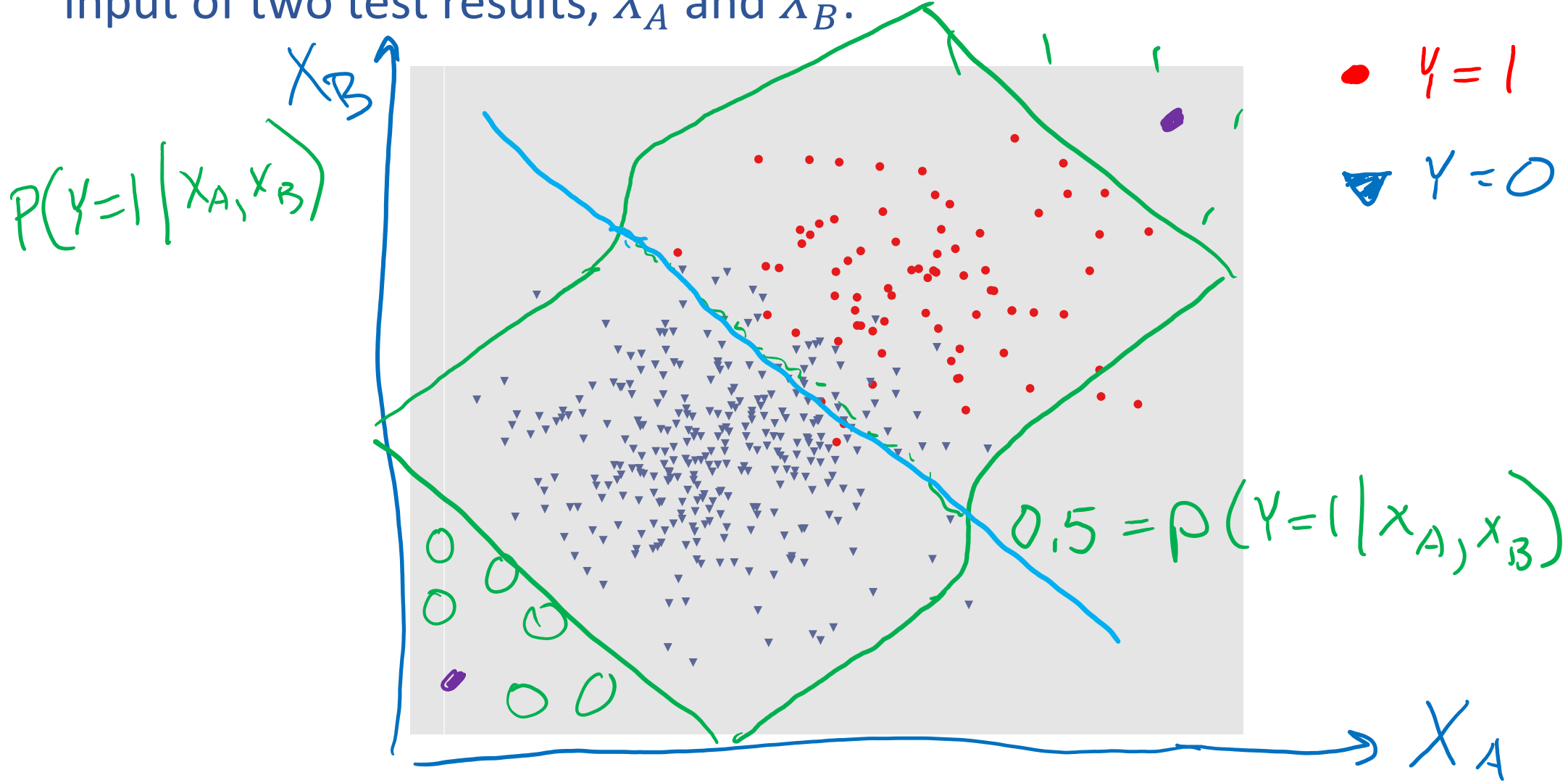
$$\ell(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{k=1}^K y_k \log \hat{y}_k$$

But now we need a model $h_{\theta}(\mathbf{x})$ that returns values that look like probabilities



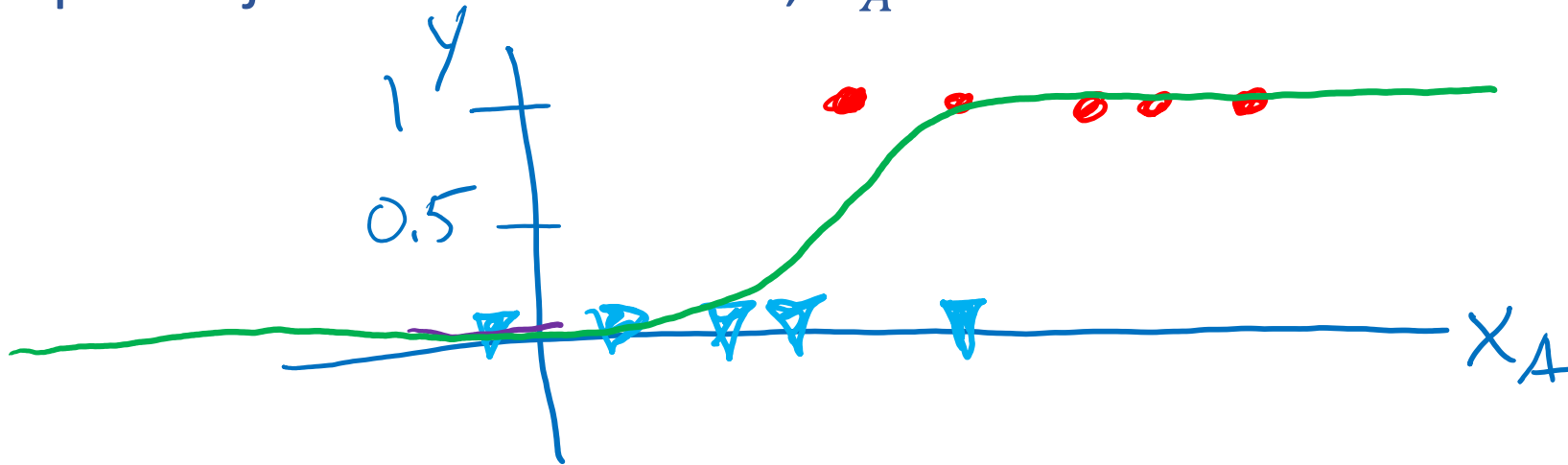
Prediction for Cancer Diagnosis

Learn to predict if a patient has cancer ($Y = 1$) or not ($Y = 0$) given the input of two test results, X_A and X_B .



Prediction for Cancer Diagnosis

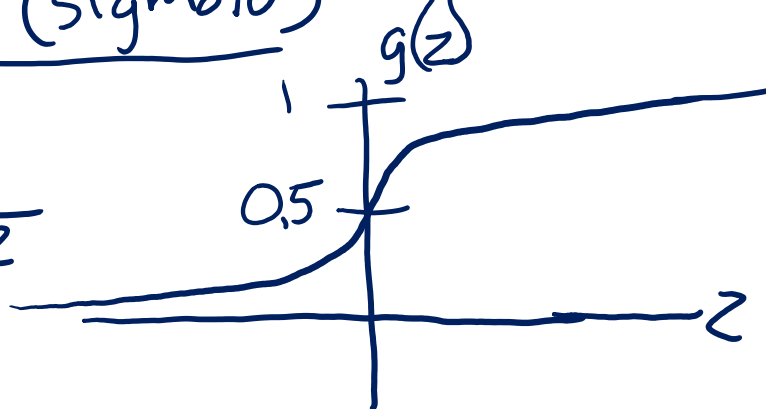
Learn to predict if a patient has cancer ($Y = 1$) or not ($Y = 0$) given the input of just one test result, X_A .



$$p(Y=1 | x_A)$$

logistic function (sigmoid)

$$g(z) = \frac{1}{1 + e^{-z}}$$

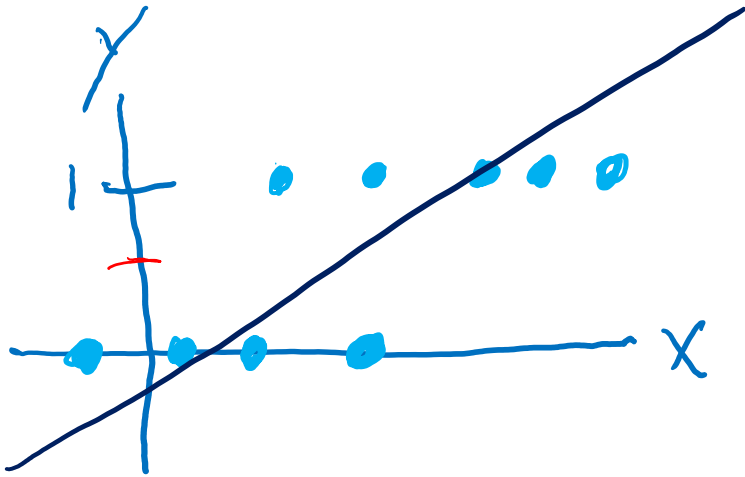


logistic regression

$$p(Y=1 | \vec{x}, \vec{\theta}) = g(\vec{\theta}^T \vec{x})$$

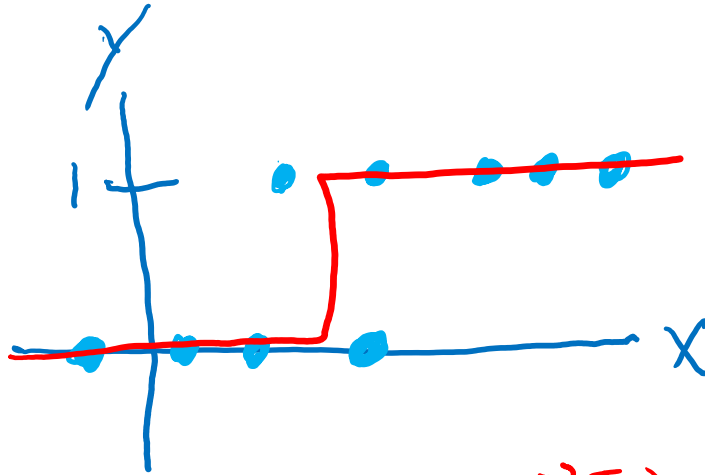
Building on a Linear Model

Linear vs Thresholded Linear vs Logistic Linear



$$\hat{y} = \vec{\theta}^T \vec{x}$$

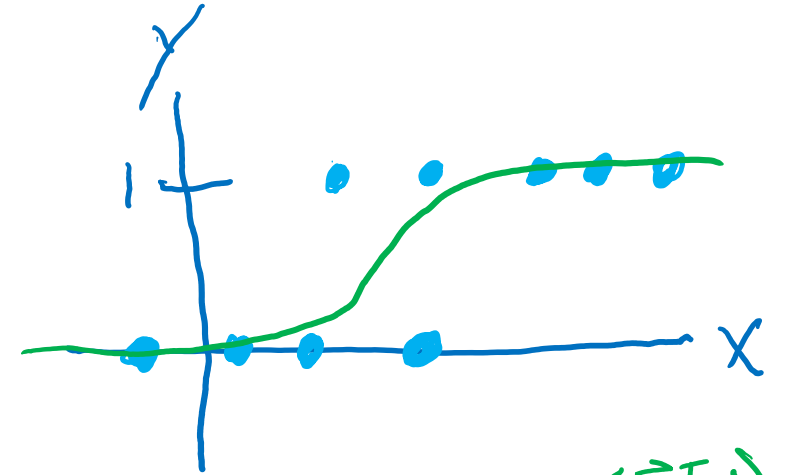
∴ not classification



$$\hat{y} = g_{\text{thresh}}(\vec{\theta}^T \vec{x})$$

∴ classification only
(0/1)

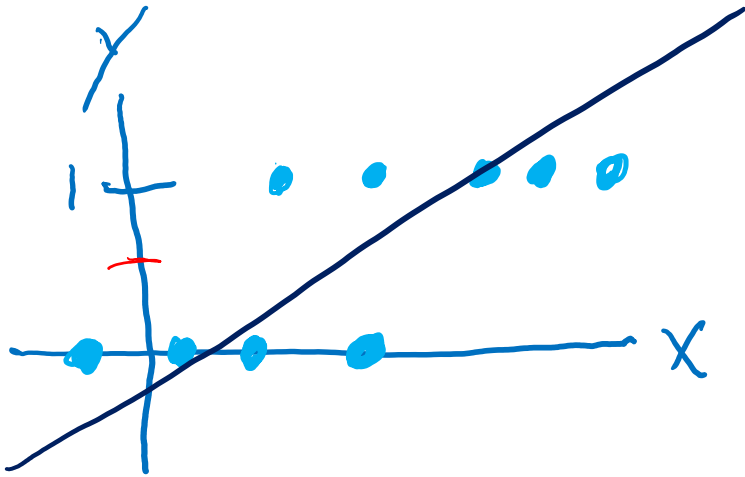
∴ zero derivatives



$$\hat{y} = g_{\text{logistic}}(\vec{\theta}^T \vec{x})$$

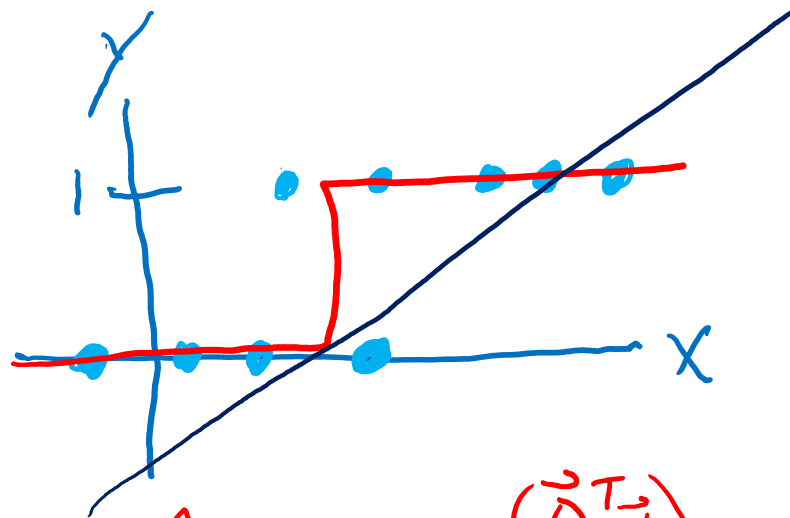
Building on a Linear Model

Linear vs Thresholded Linear vs Logistic Linear



$$\hat{y} = \vec{\theta}^T \mathbf{x}$$

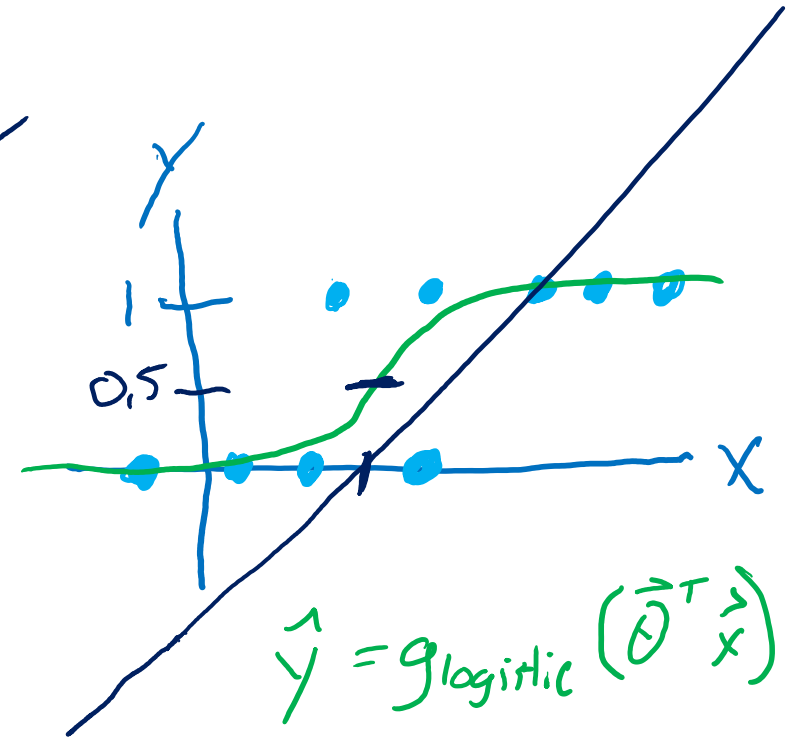
∴ not classification



$$\hat{y} = g_{\text{thresh}}(\vec{\theta}^T \mathbf{x})$$

∴ classification only
(0/1)

∴ zero derivatives



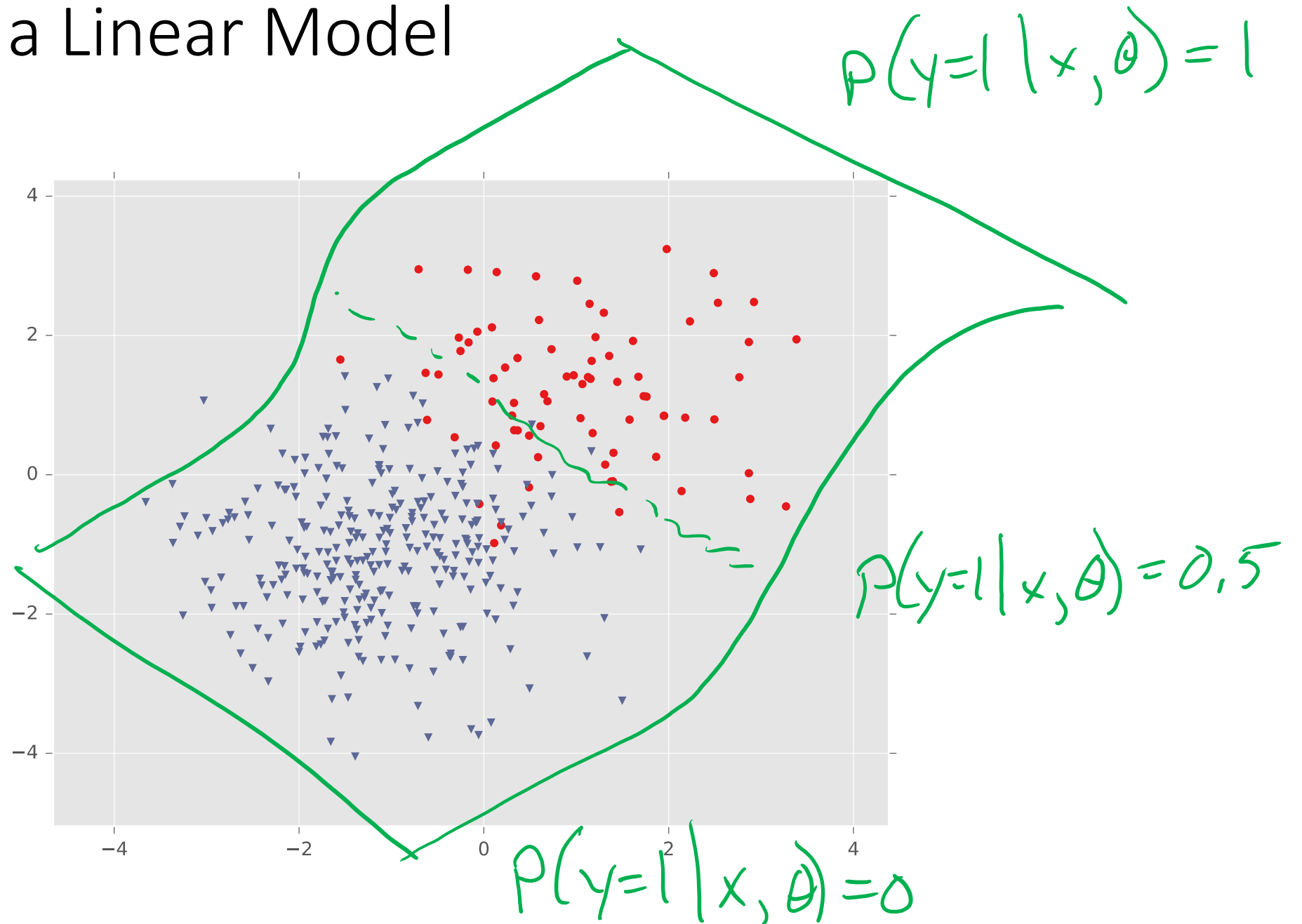
$$\hat{y} = g_{\text{logistic}}(\vec{\theta}^T \mathbf{x})$$

$$\uparrow \sigma(\vec{\theta}^T \mathbf{x})$$

Building on a Linear Model

$$\vec{x} = \begin{bmatrix} 1 \\ x_A \\ x_B \end{bmatrix}$$

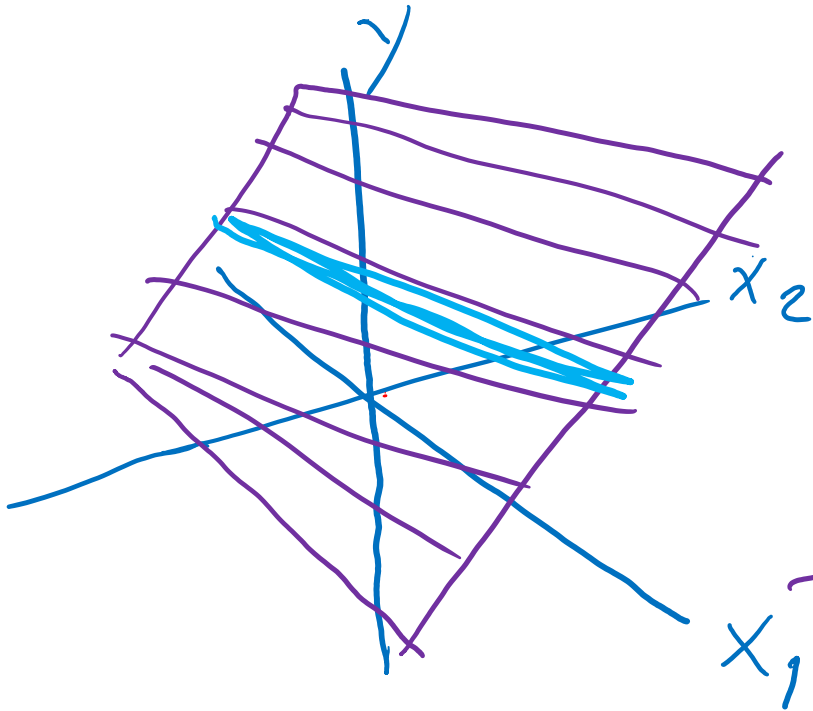
$$\vec{\theta} = \begin{bmatrix} b \\ w_1 \\ w_2 \end{bmatrix}$$



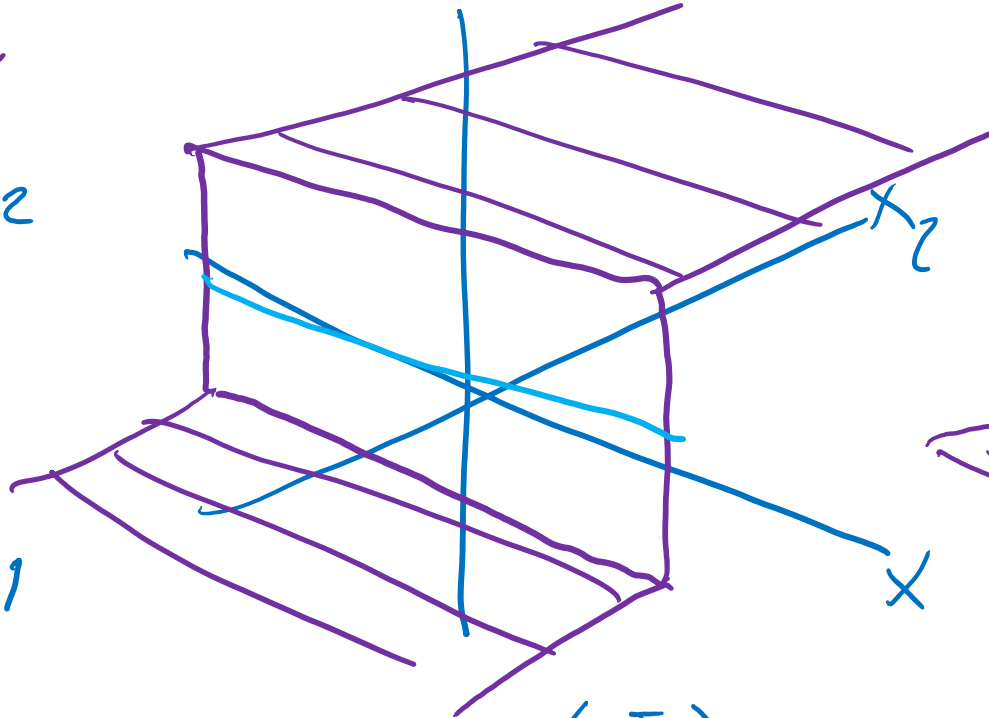
Building on a Linear Model

$$\frac{d\sigma}{dz} = \sigma(z)(1 - \sigma(z))$$

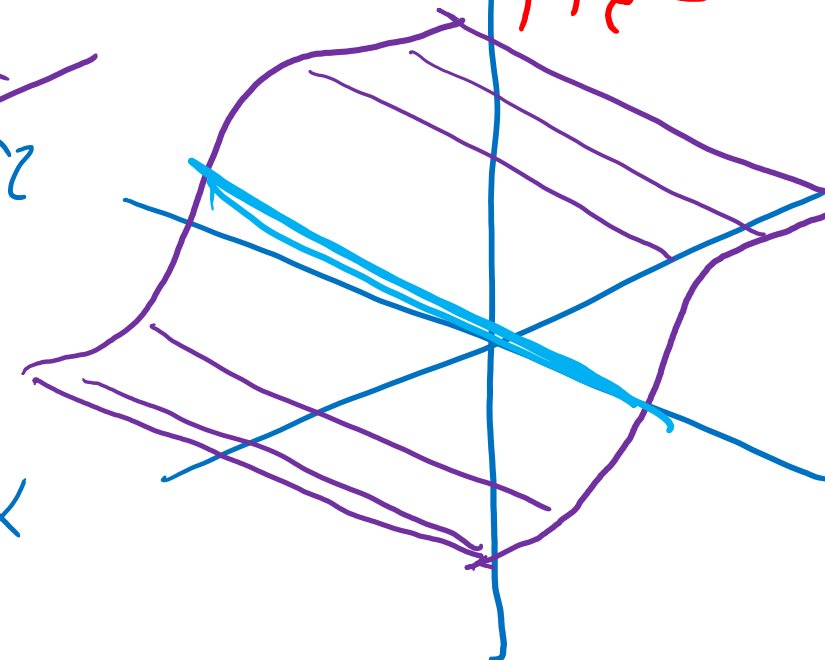
$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



$$y = \vec{\theta}^T \vec{x}$$



$$y = \text{sign}(\theta^T x)$$



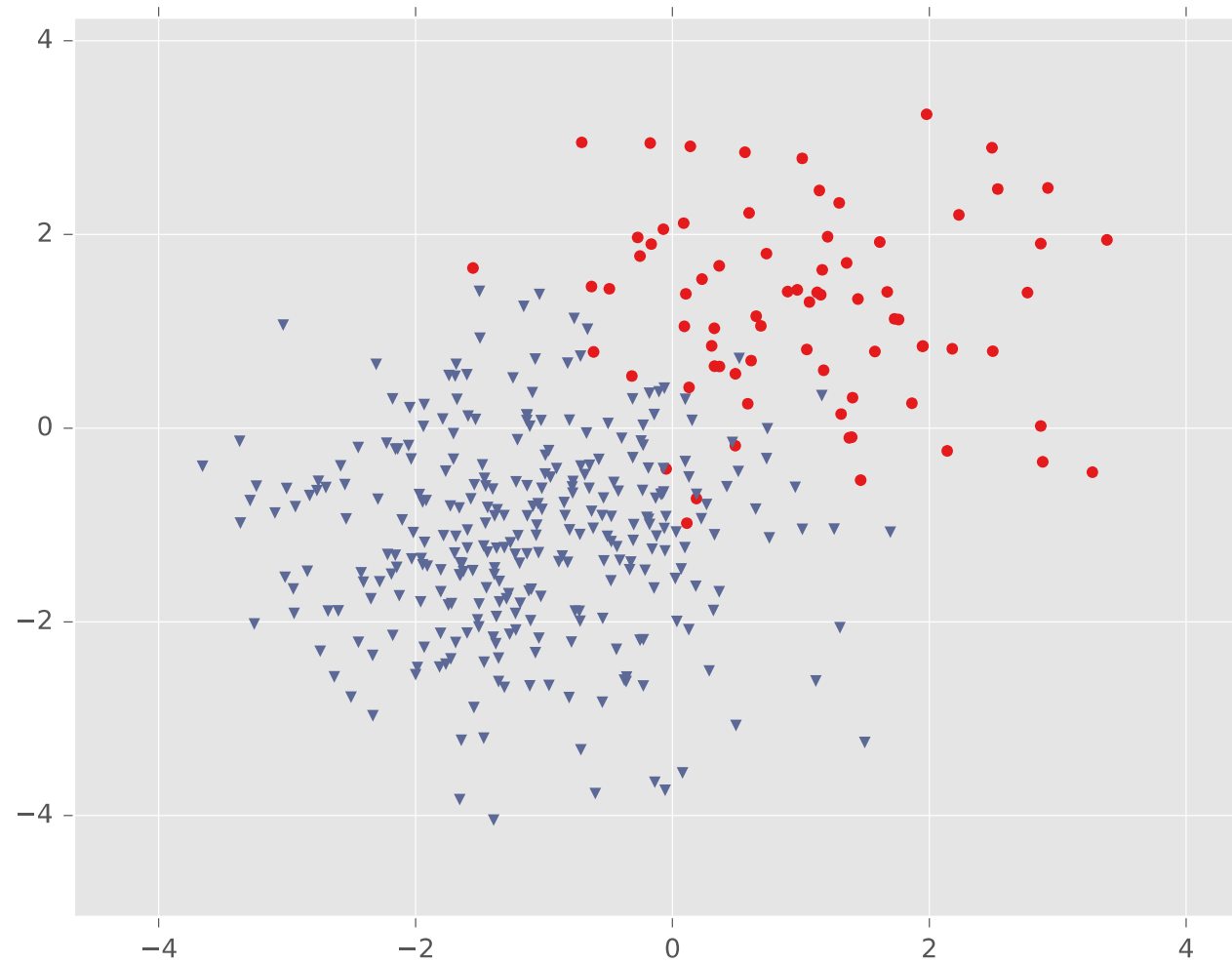
$$y = \sigma(\theta^T x)$$

logistic
function

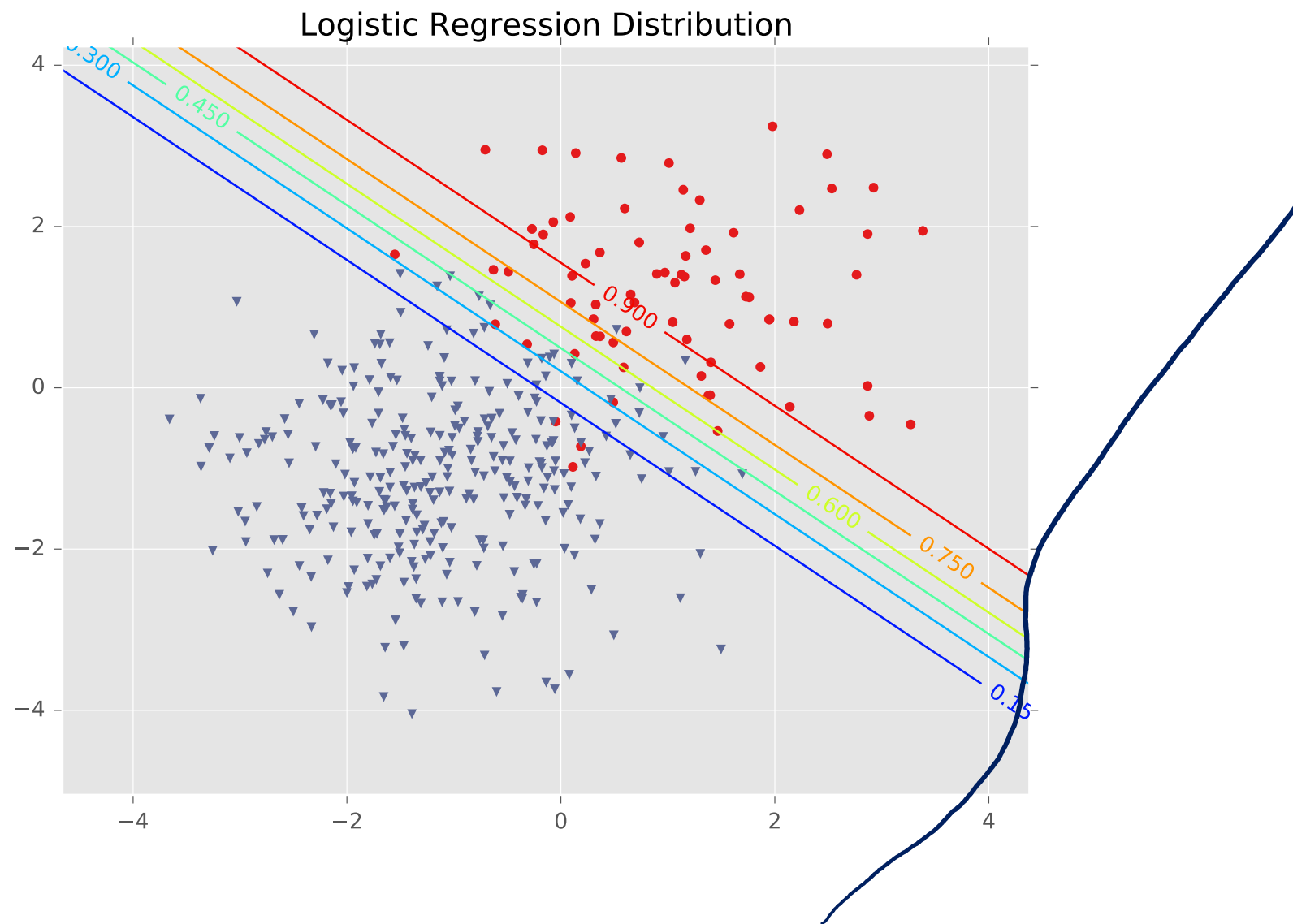
$$p(y|x)$$

$y=1$

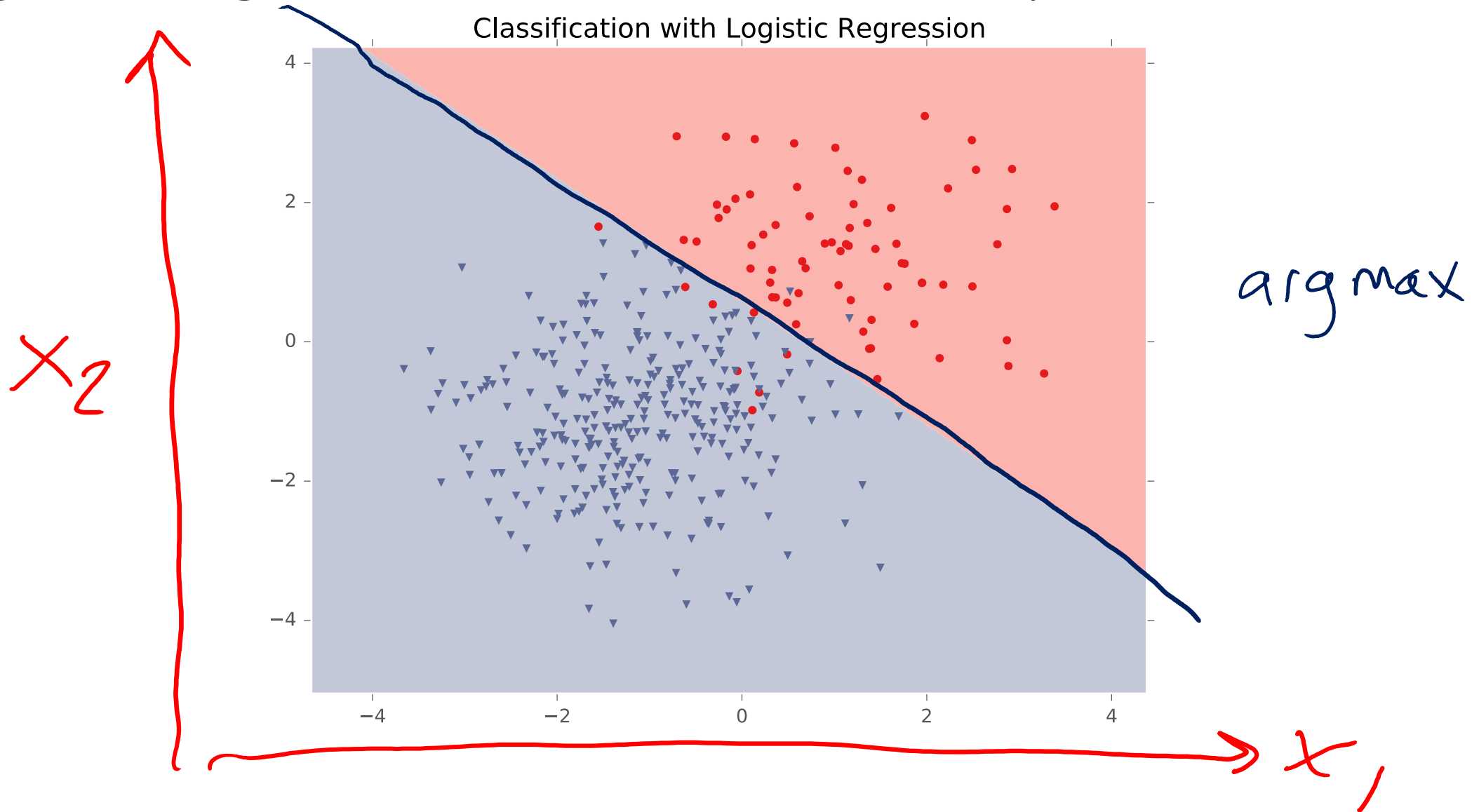
Logistic Regression



Logistic Regression



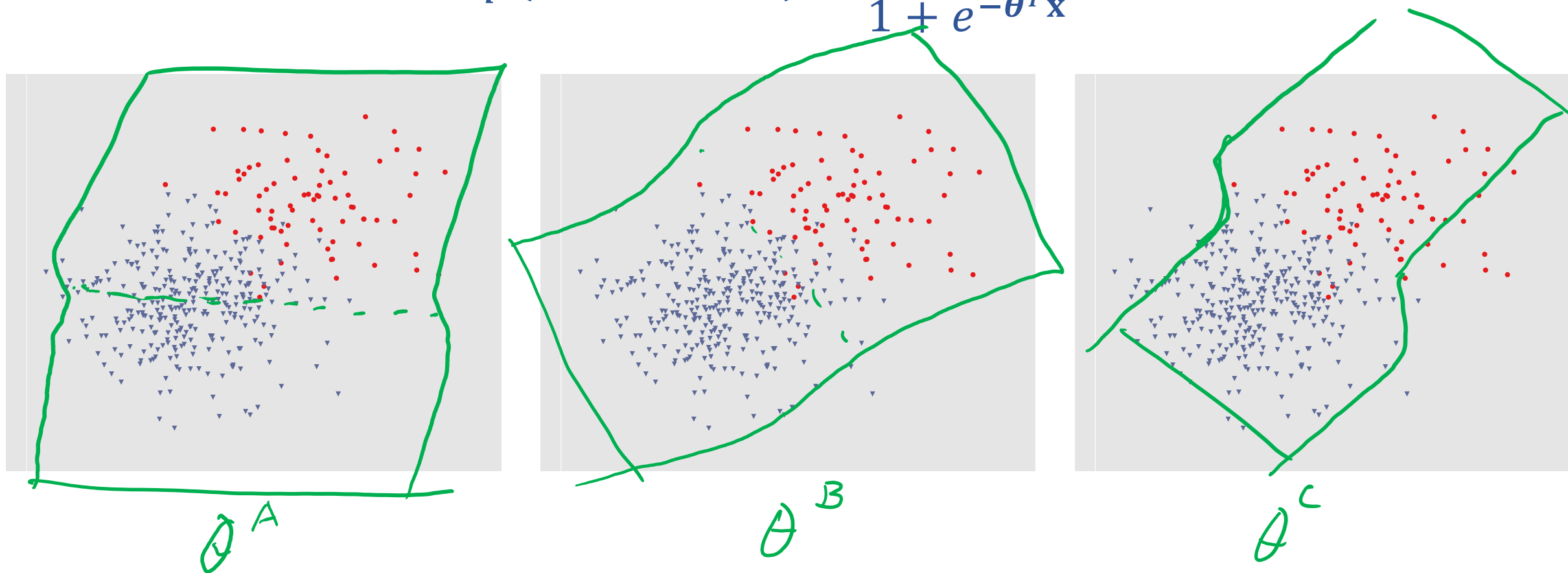
Logistic Regression Decision Boundary



Optimizing a Model for Cancer Diagnosis

Learn to predict if a patient has cancer ($Y = 1$) or not ($Y = 0$) given the input of two test results, X_A , X_B . Note: bias term included in \mathbf{x} .

$$p(Y = 1 | \mathbf{x}, \boldsymbol{\theta}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^T \mathbf{x}}}$$



Binary Logistic Regression

$$\hat{y} \in [0, 1]$$

1) Model

$$\hat{y} = h(\vec{x}) = P(Y=1|\vec{x}) = \frac{1}{1 + e^{-\theta^T \vec{x}}}$$
$$P(Y=0|\vec{x}) = 1 - \left[\frac{1}{1 + e^{-\theta^T \vec{x}}} \right]$$

2) Objective function

$$\frac{1}{N} \sum_{i=1}^N \ell(y^{(i)}, \hat{y}^{(i)}) = - \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K y_k^{(i)} \log \hat{y}_k$$
$$= - \frac{1}{N} \sum \left[y_1^{(i)} \log \hat{y}_1 + y_2^{(i)} \log \hat{y}_2 \right]$$
$$= - \frac{1}{N} \sum \left[y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log (1 - \hat{y}^{(i)}) \right]$$

3) Solve for $\hat{\theta}$

Binary Logistic Regression

Gradient

$$J(\theta) = -\frac{1}{N} \sum_{i=1}^N \log p(y^{(i)} | \vec{x}^{(i)}, \vec{\theta})$$

$$\nabla_{\theta} J^{(i)} = \begin{bmatrix} \frac{\partial J^{(i)}}{\partial \theta_1} \\ \frac{\partial J^{(i)}}{\partial \theta_m} \\ \vdots \end{bmatrix}$$

SGD

$$J^{(i)}(\vec{\theta})$$

$$\frac{\partial J^{(i)}}{\partial \theta_m}$$

$$= - (y^{(i)} - \sigma(\theta^T x^{(i)})) x_m^{(i)}$$

$$\nabla_{\theta} J^{(i)} = - (y^{(i)} - \sigma(\theta^T x^{(i)})) \vec{x}^{(i)}$$

$$\vec{x} \in \mathbb{R}^m \quad \vec{\theta} \in \mathbb{R}^m$$

$$\nabla_{\theta} J^{(i)}(\theta):$$


$$\mathbb{R}^m \rightarrow \mathbb{R}^m$$

$$\begin{bmatrix} x_1^{(i)} \\ \vdots \\ x_m^{(i)} \end{bmatrix}$$

Solve Logistic Regression

$J(\theta)$ is convex
😊

$$\hat{y} = g(\boldsymbol{\theta}^T \mathbf{x}) \quad g(z) = \frac{1}{1+e^{-z}}$$

$$J(\boldsymbol{\theta}) = -\frac{1}{N} \sum_i (y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}))$$


$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = -\frac{1}{N} \sum_i (y^{(i)} - \hat{y}^{(i)}) \mathbf{x}^{(i)}$$

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = 0?$$

No closed form solution ☹️



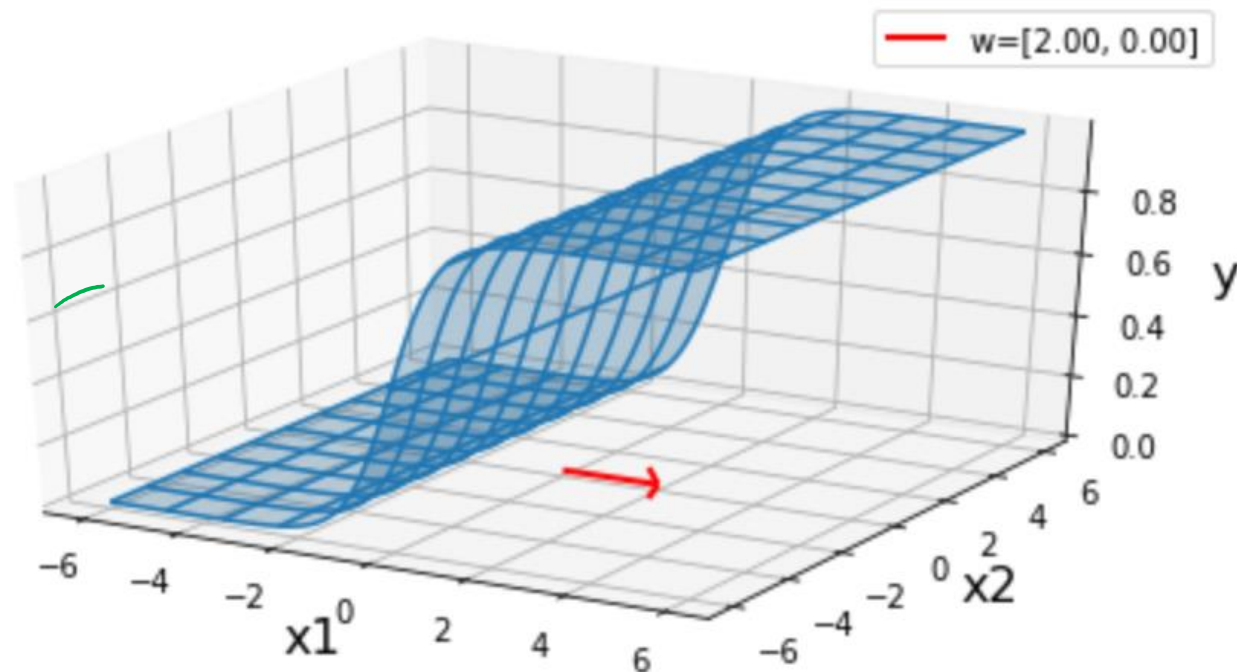
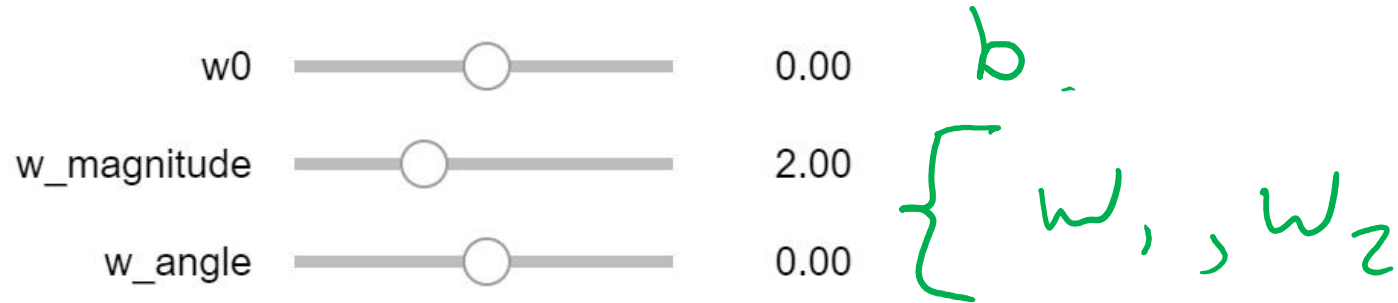
Back to iterative methods. Solve with (stochastic) gradient descent, Newton's method, or Iteratively Reweighted Least Squares (IRLS)

Logistic Regression Decision Boundary



Exercise

Interact with the linear_logistic.ipynb posted on the course website schedule



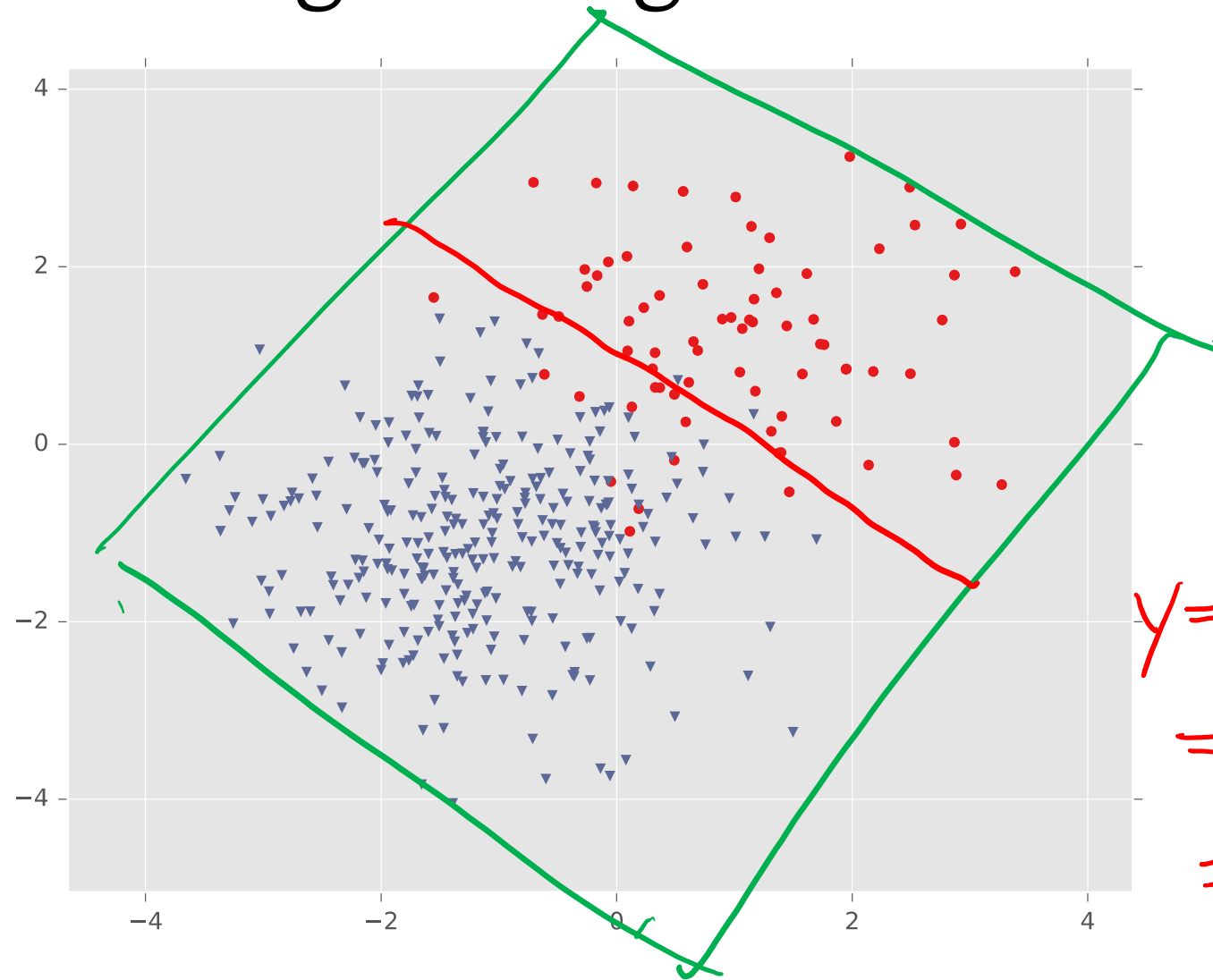
Linear in Higher Dimensions

1-D $y = w x + b$
2-D $y = w_1 x_1 + w_2 x_2 + b$

What are these linear shapes called for 1-D, 2-D, 3-D, M-D input?

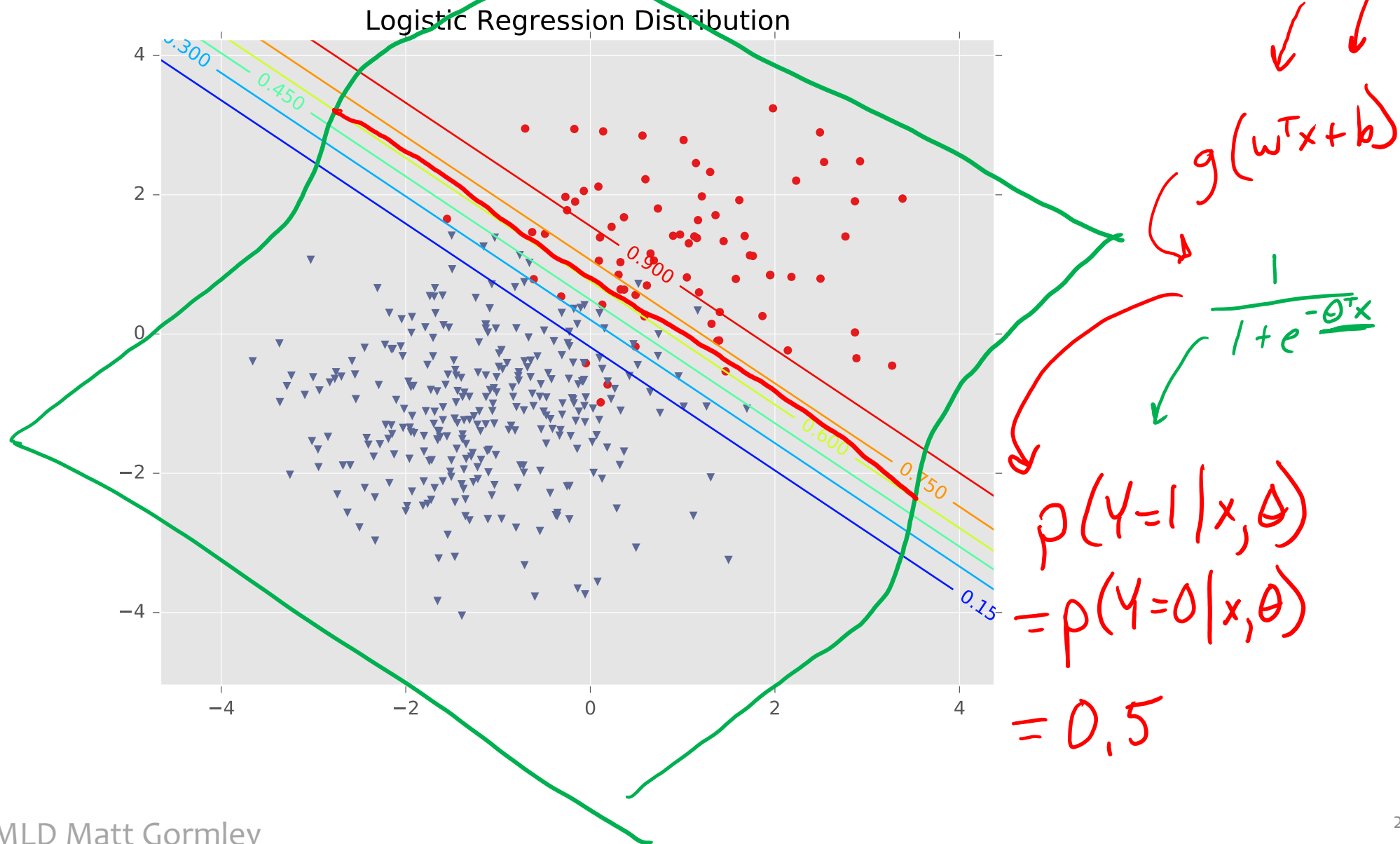
	$x \in \mathbb{R}$	$x \in \mathbb{R}^2$	$x \in \mathbb{R}^3$	$x \in \mathbb{R}^M$
$\rightarrow y = w^T x + b$	line	plane	hyperplane	hyperplane
$w^T x + b = 0$	point	line	plane	hyperplane
$w^T x + b \geq 0$	halfline	halfplane	halfspace	halfspace

Logistic Regression



$$\begin{aligned} y &= \theta^T x \\ &= w^T x + b \\ &= 0 \end{aligned}$$

Logistic Regression



Logistic Regression

Classification with Logistic Regression

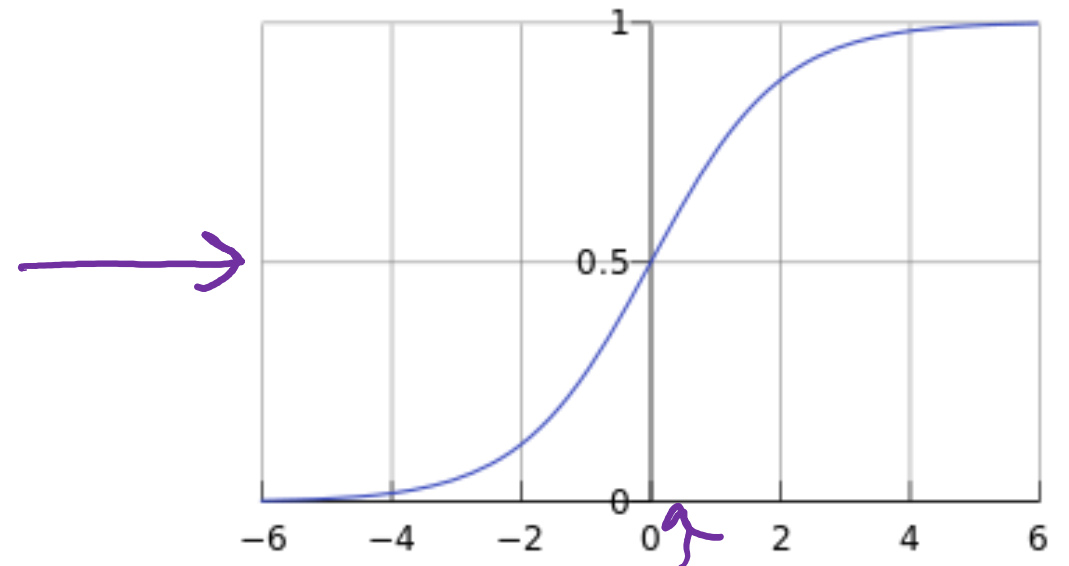


Poll 1

For a point \mathbf{x} on the decision boundary of logistic regression,

does $\underbrace{g(\mathbf{w}^T \mathbf{x} + b)}_{0.5} = \underbrace{\mathbf{w}^T \mathbf{x} + b}_0$?

\downarrow
 $g(z) = \frac{1}{1 + e^{-z}}$

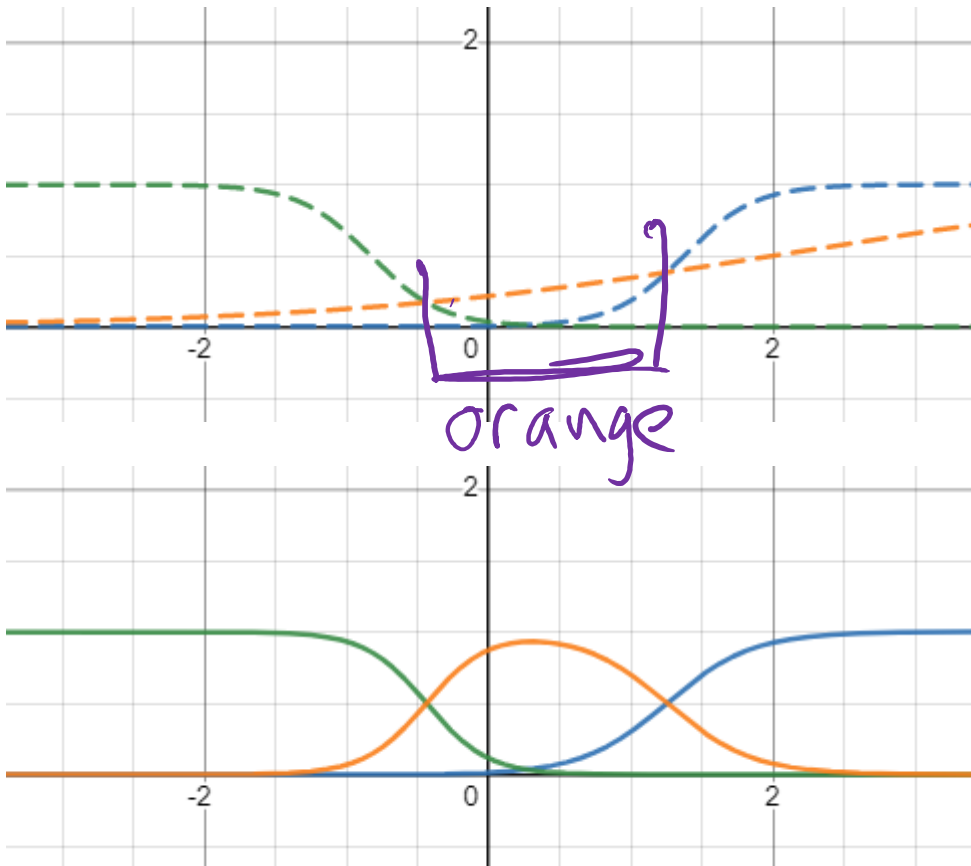


Multi-class Logistic Regression

Multi-class Logistic Regression

Desmos Demo:

<https://www.desmos.com/calculator/53bautbxjp>



$$\begin{aligned} g_0(\theta_1^T x) & \leftarrow \\ g_1(\theta_2^T x) & \leftarrow \\ g_2(\theta_3^T x) & \leftarrow \end{aligned}$$

$$\hat{y} =$$

$$g_{\text{sm}} \left(\begin{bmatrix} \theta_1^T x \\ \theta_2^T x \\ \theta_3^T x \end{bmatrix} \right)$$

Multi-class Logistic Regression

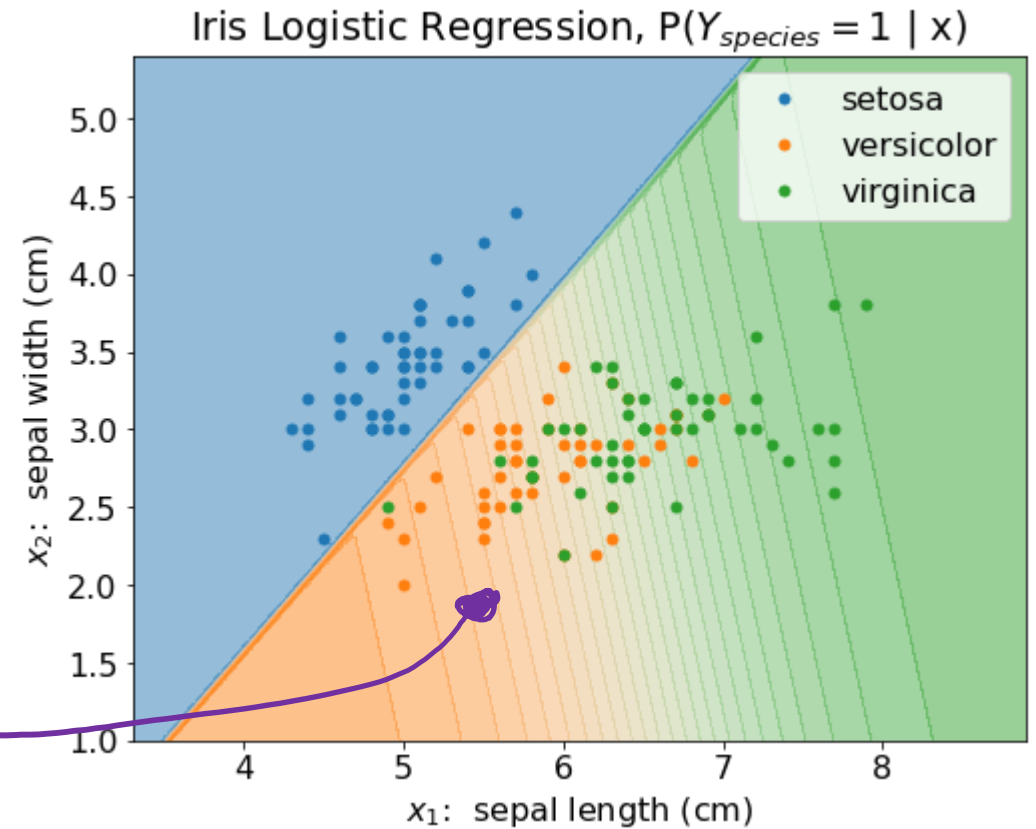
Cross-entropy loss

$$\ell(\mathbf{y}, \hat{\mathbf{y}}) = -\sum_{k=1}^K y_k \log \hat{y}_k$$

Model

$$\hat{\mathbf{y}} = g_{\mathbf{w}, \mathbf{b}}$$

$$\begin{bmatrix} 0 \\ .8 \\ .2 \end{bmatrix}$$



Logistic Function

Logistic (sigmoid) function converts value from $(-\infty, \infty) \rightarrow (0, 1)$

$$g(z) = \frac{1}{1 + e^{-z}} = \frac{e^z}{e^z + 1}$$

$g(z)$ and $1 - g(z)$ sum to one

Example 2 $\rightarrow g(2) = 0.88, \quad 1 - g(2) = 0.12$

Softmax Function

Softmax function convert each value in a vector of values from $(-\infty, \infty) \rightarrow (0, 1)$, such that they all sum to one.

$$g(z)_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}}$$

$$\begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_K \end{bmatrix} \rightarrow \begin{bmatrix} e^{z_1} \\ e^{z_2} \\ \vdots \\ e^{z_K} \end{bmatrix} \cdot \frac{1}{\sum_{k=1}^K e^{z_k}}$$

Example

$$\begin{bmatrix} -1 \\ 4 \\ 1 \\ -2 \\ 3 \end{bmatrix} \rightarrow \begin{bmatrix} 0.0047 \\ 0.7008 \\ 0.0349 \\ 0.0017 \\ 0.2578 \end{bmatrix}$$

Multiclass Predicted Probability

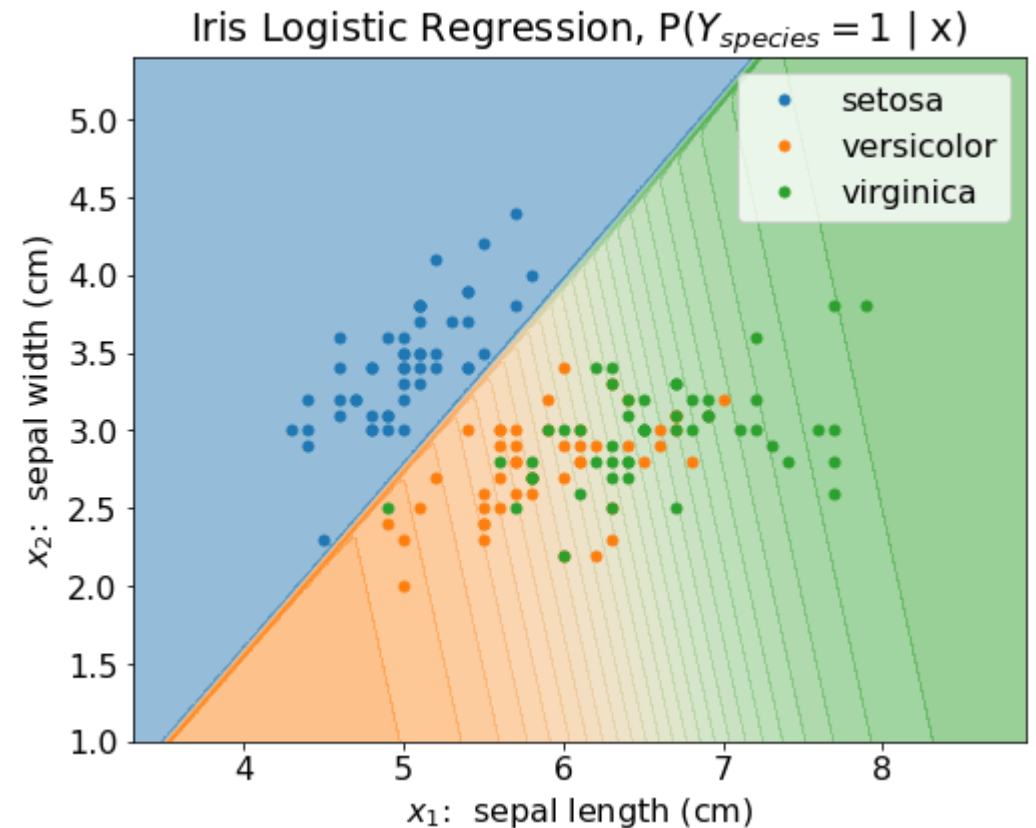
Multiclass logistic regression uses the parameters learned across all K classes to predict the discrete conditional probability distribution of the output Y given a specific input vector \mathbf{x}

$$\begin{bmatrix} p(Y = 1 \mid \mathbf{x}, \boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3) \\ p(Y = 2 \mid \mathbf{x}, \boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3) \\ p(Y = 3 \mid \mathbf{x}, \boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3) \end{bmatrix} = \begin{bmatrix} e^{\boldsymbol{\theta}_1^T \mathbf{x}} \\ e^{\boldsymbol{\theta}_2^T \mathbf{x}} \\ e^{\boldsymbol{\theta}_3^T \mathbf{x}} \end{bmatrix} \cdot \frac{1}{\sum_{k=1}^K e^{\boldsymbol{\theta}_k^T \mathbf{x}}}$$

each of these
is a vector (confusing
notation)

Multiclass Predicted Probability

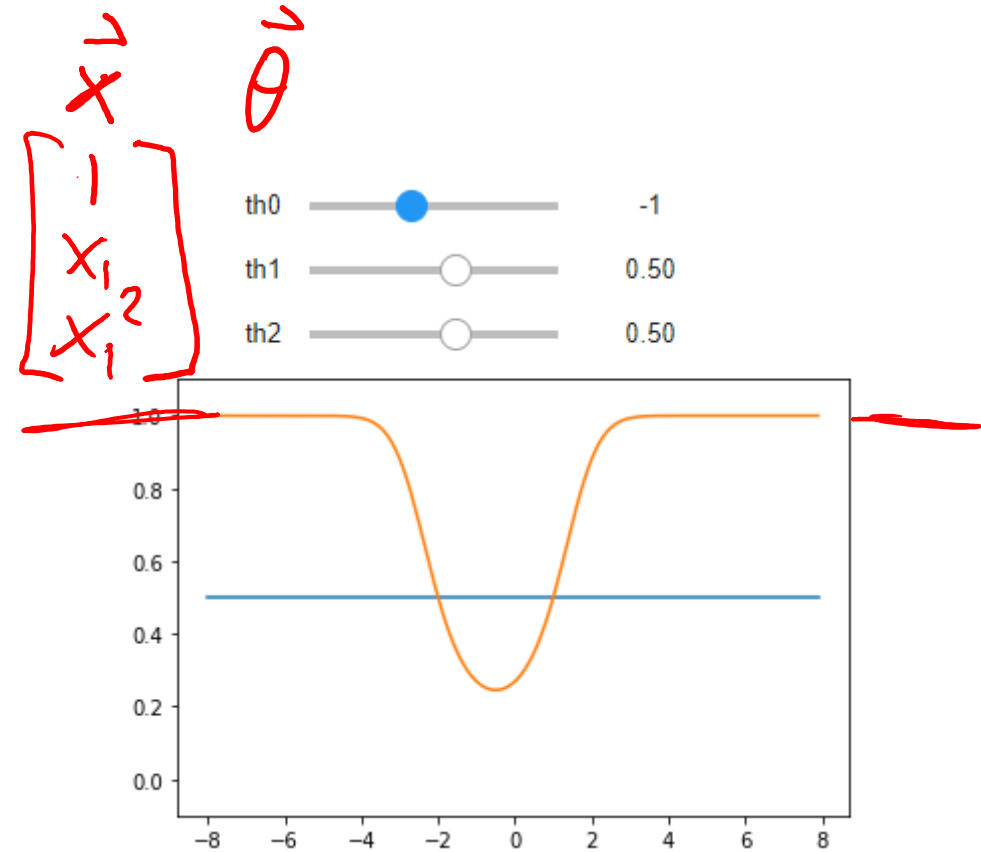
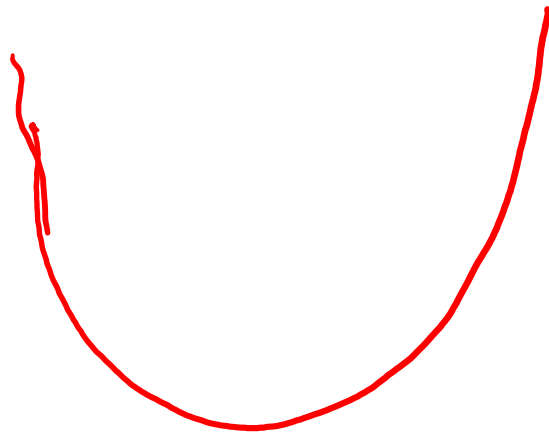
Multiclass logistic regression uses the parameters learned across all K classes to predict the discrete conditional probability distribution of the output Y given a specific input vector \mathbf{x}



Logistic Regression with Polynomial Features

Exercise

Interact with the `logistic_quadratic.ipynb` posted on the course website schedule



Exercise

Interact with the `logistic_quadratic.ipynb` posted on the course website schedule

