

10-315 Introduction to ML

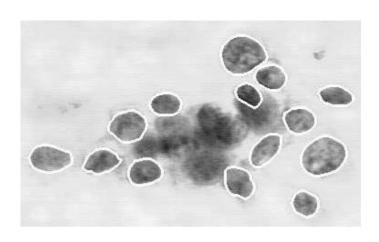
Logistic Regression

Instructor: Pat Virtue

Example: Breast cancer classification

Well-known classification example: using machine learning to diagnose whether a breast tumor is benign or malignant [Street et al., 1992]

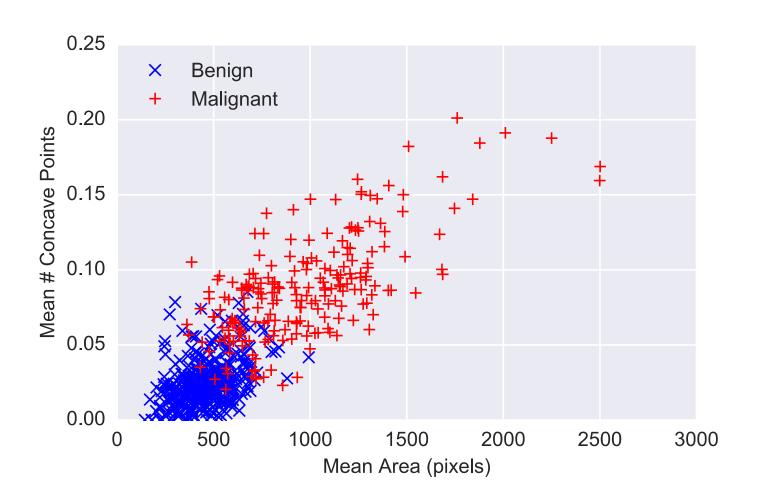
Setting: doctor extracts a sample of fluid from tumor, stains cells, then outlines several of the cells (image processing refines outline)



System computes features for each cell such as area, perimeter, concavity, texture (10 total); computes mean/std/max for all features

Example: Breast cancer classification

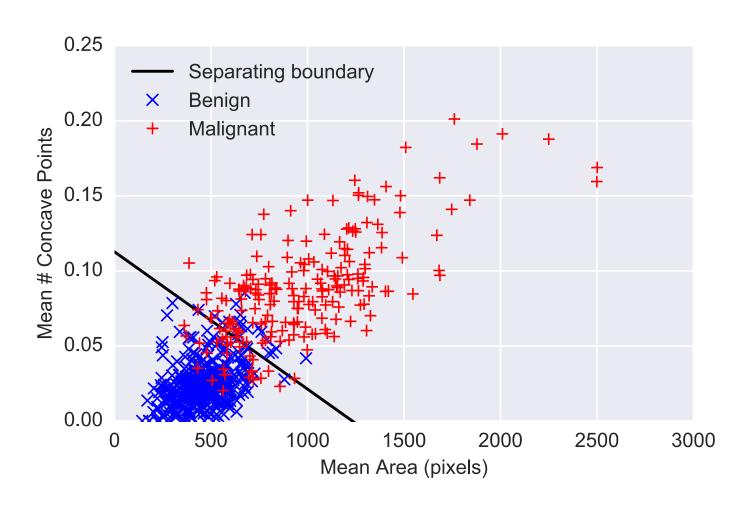
Plot of two features: mean area vs. mean concave points, for two classes



Slide credit: CMU Al Zico Kolter

Linear classification example

Linear classification: linear decision boundary

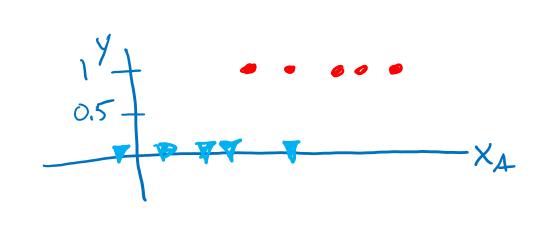


Logistic regression for classification

Linear classification: linear decision boundary

Probabilistic classification: provide $P(Y = 1 \mid x)$ rather than just $\hat{y} \in \{0, 1\}$



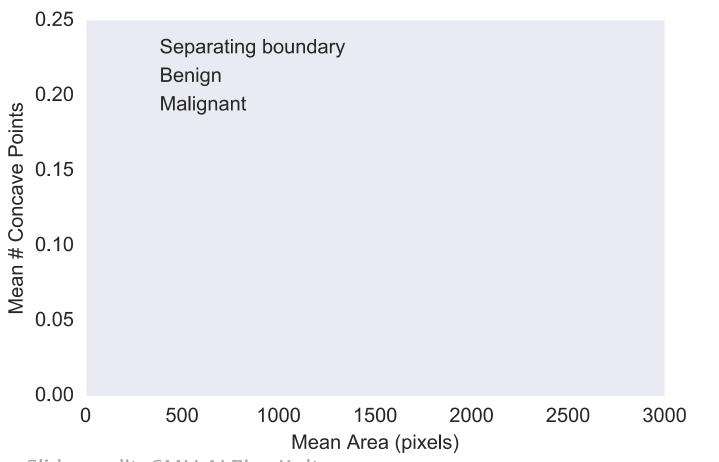


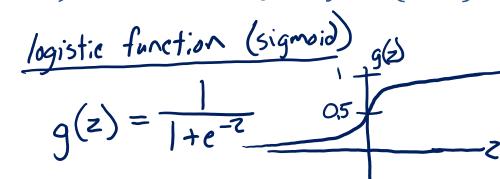
Slide credit: CMU AI Zico Kolter

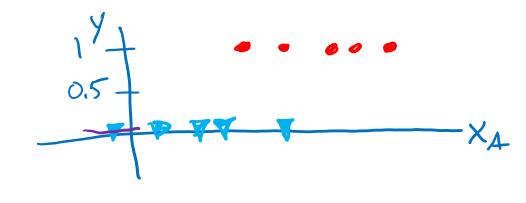
Logistic regression for classification

Linear classification: linear decision boundary

Probabilistic classification: provide $P(Y = 1 \mid x)$ rather than just $\hat{y} \in \{0, 1\}$





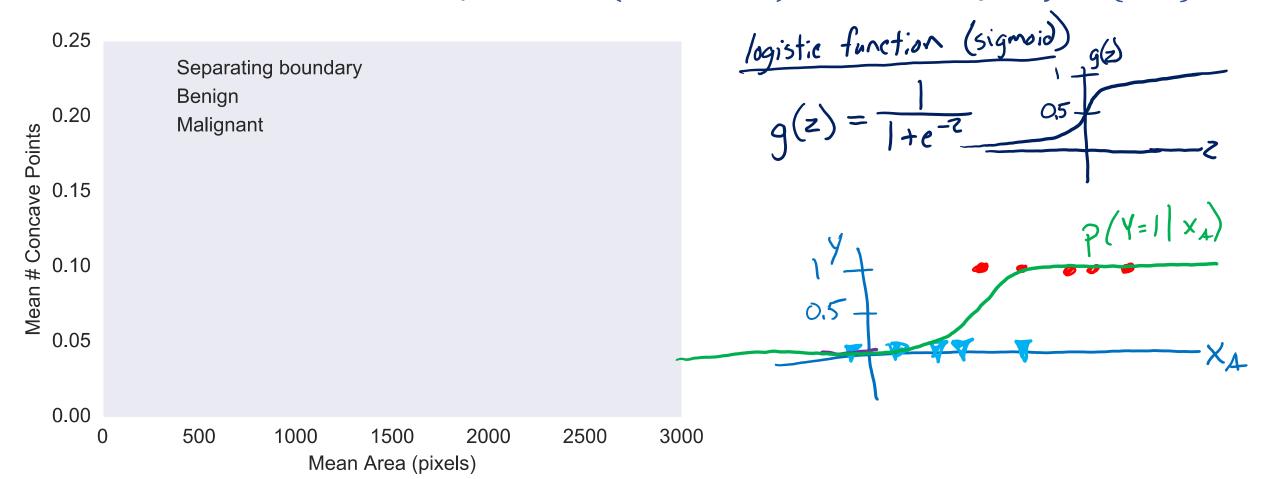


Slide credit: CMU AI Zico Kolter

Logistic regression for classification

Linear classification: linear decision boundary

Probabilistic classification: provide $P(Y = 1 \mid x)$ rather than just $\hat{y} \in \{0, 1\}$



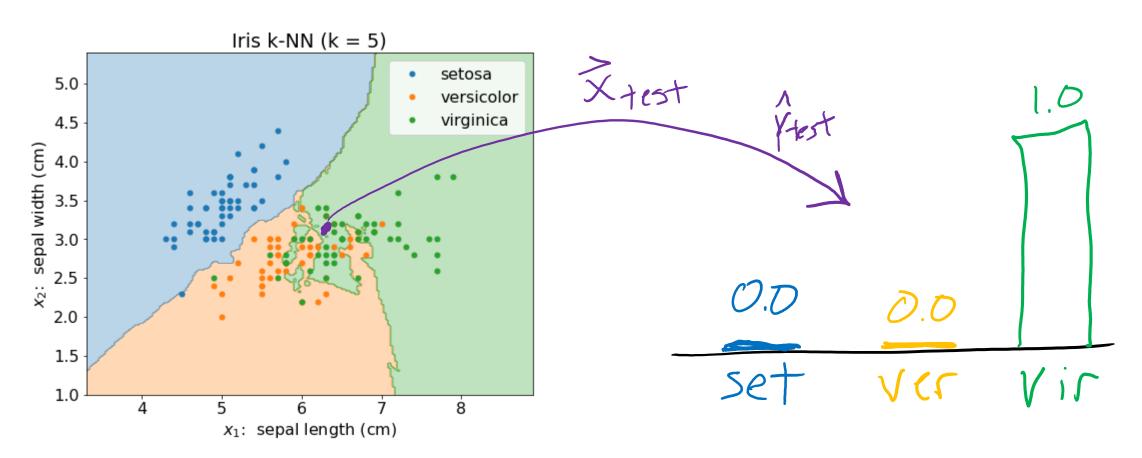
Slide credit: CMU AI Zico Kolter

Pre-reading

Cross-entropy loss

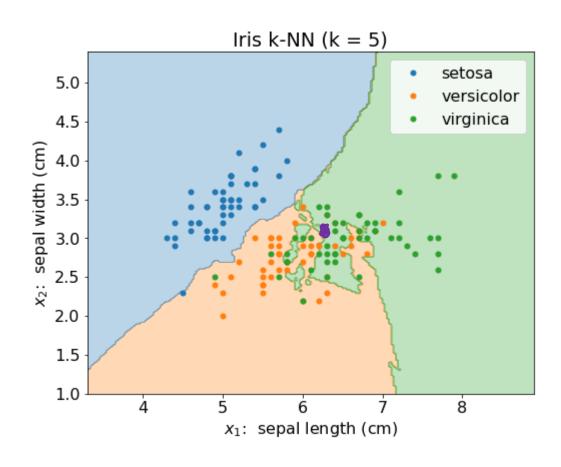
Classification Decisions

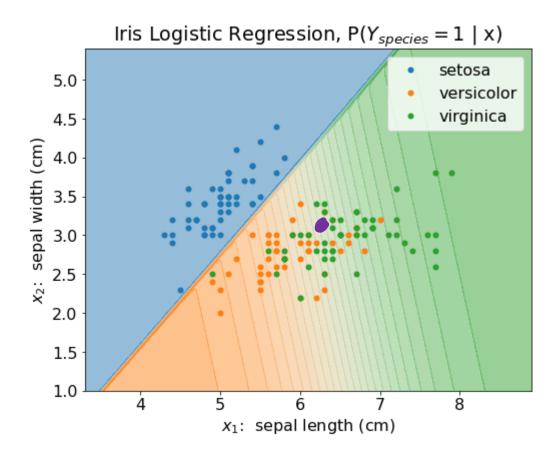
Predicting one specific class is troubling, especially when we know that there is some uncertainty in our prediction



Classification Probability

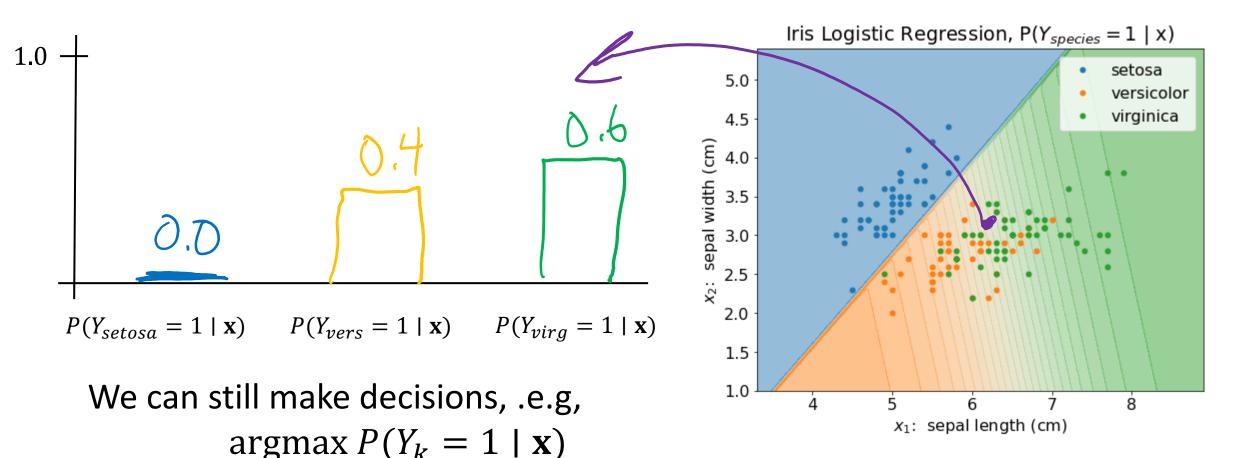
Constructing a model than can return the probability of the output being a specific class could be incredibly useful





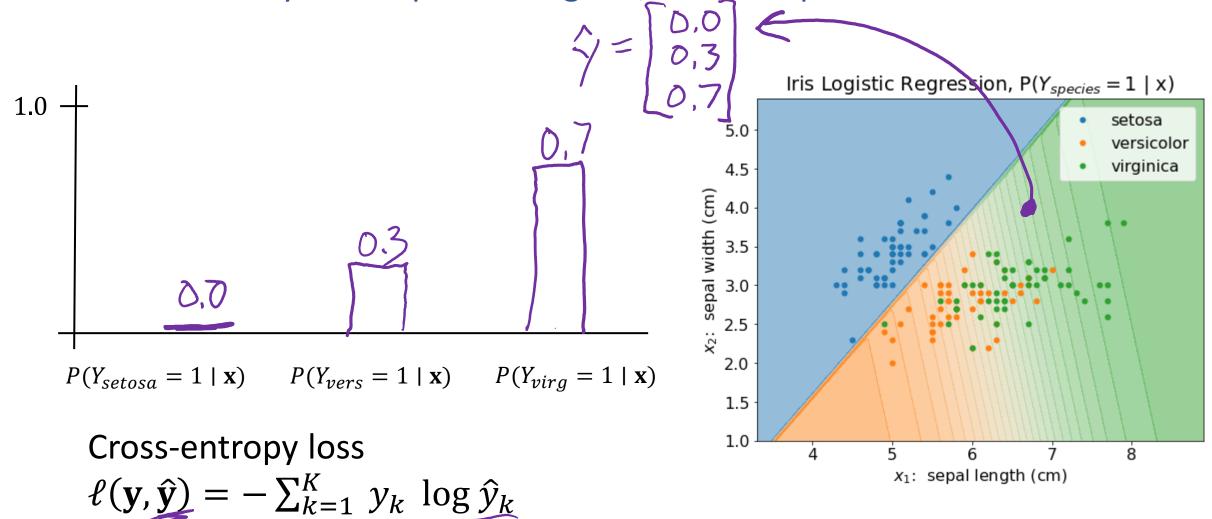
Classification Probability

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



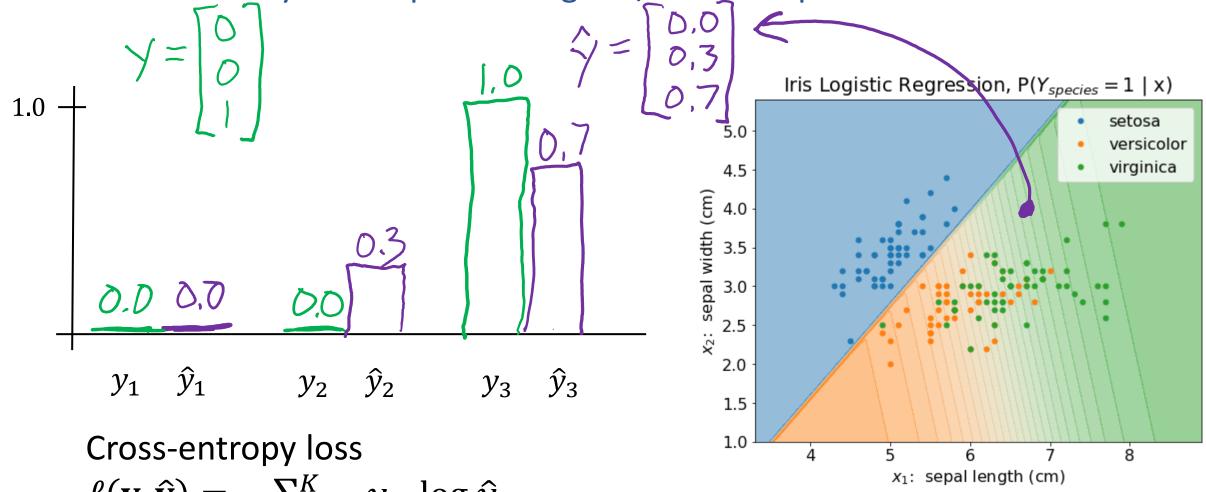
Loss for Probabilty Disributions

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



Loss for Probabilty Disributions

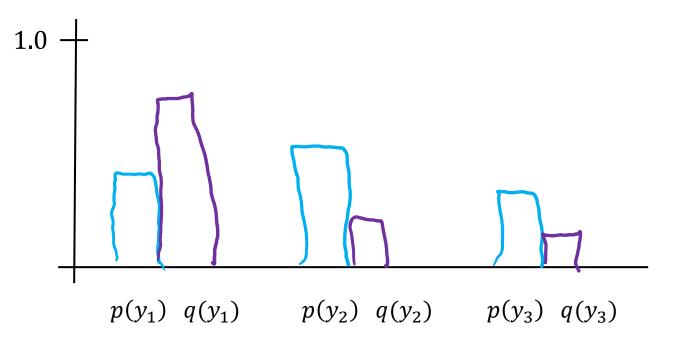
We need a way to compare how good/bad each prediction is



Loss for Probabilty Disributions

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

distributions*



*when used in logistic regression **y** is always a one-hot vector

Cross-entropy loss

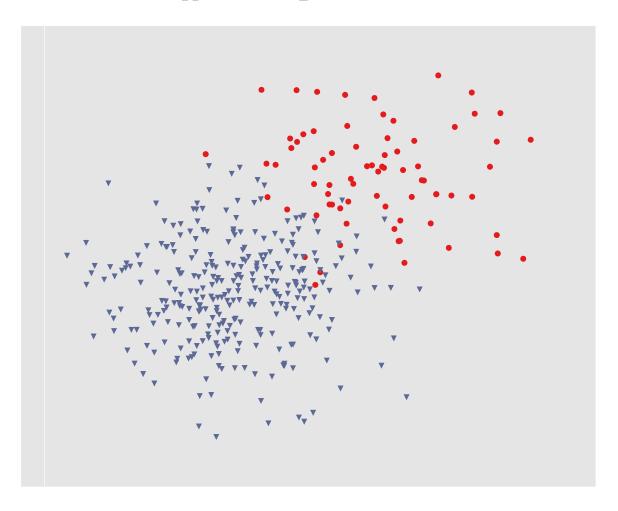
$$H(P,Q) = -\sum_{k=1}^{K} p(y_k) \log q(y_k)$$

Pre-reading

Linear model for classification

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 .

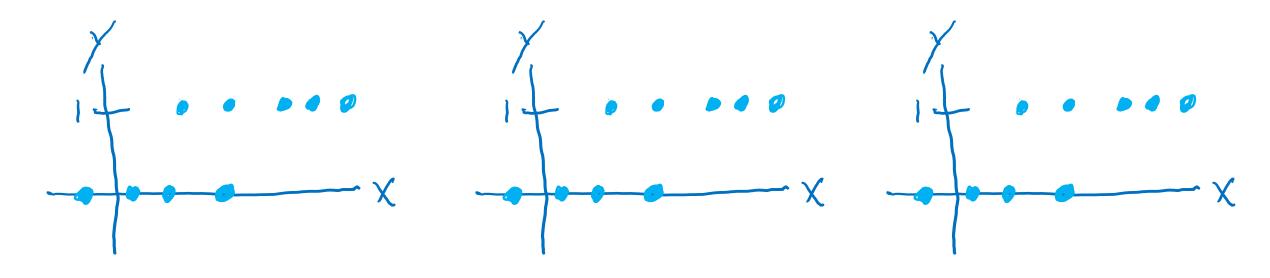


Prediction for Cancer Diagnosis

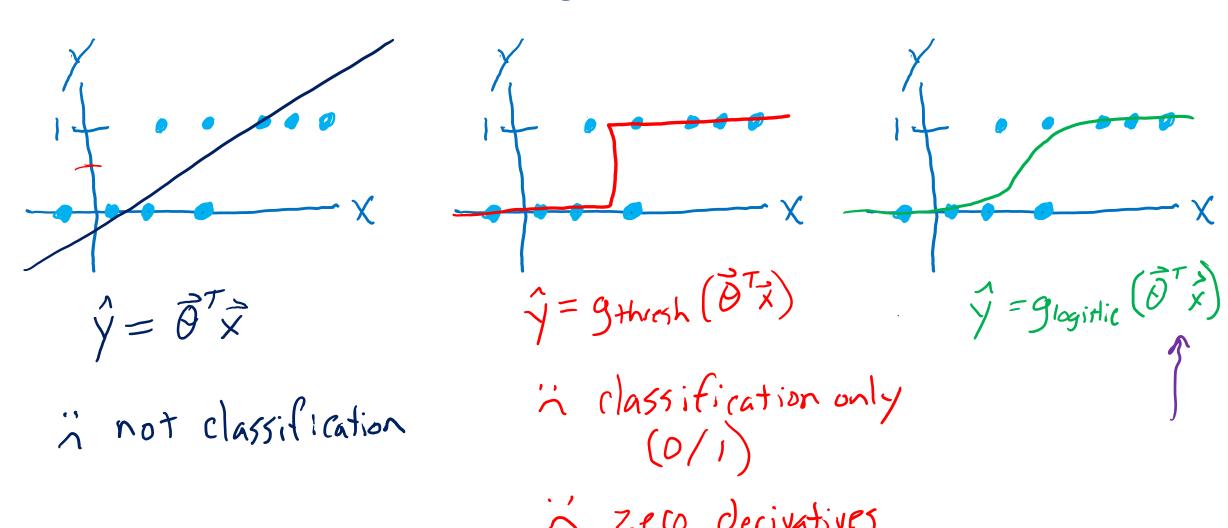
Learn to predict if a patient has cancer (Y = 1) or not (Y = 0) given the input of just one test result, X_A .

$$\frac{\log istic \ function \ (sigmoid)}{q(z) = \frac{1}{1 + e^{-2}}} \xrightarrow{0.5} \frac{\log istic \ legression}{q(\vec{\sigma} \times \vec{\sigma}) = q(\vec{\sigma} \times \vec{\sigma})}$$

Linear vs Thresholded Linear vs Logistic Linear



Linear vs Thresholded Linear vs Logistic Linear

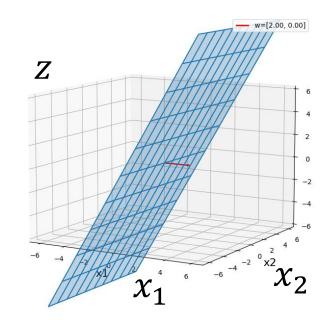


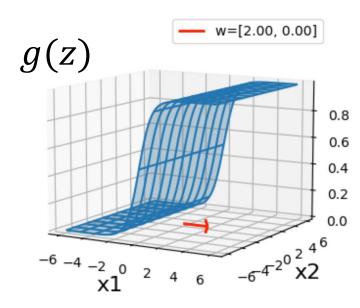
Linear model for classification (now with 2+ input features)

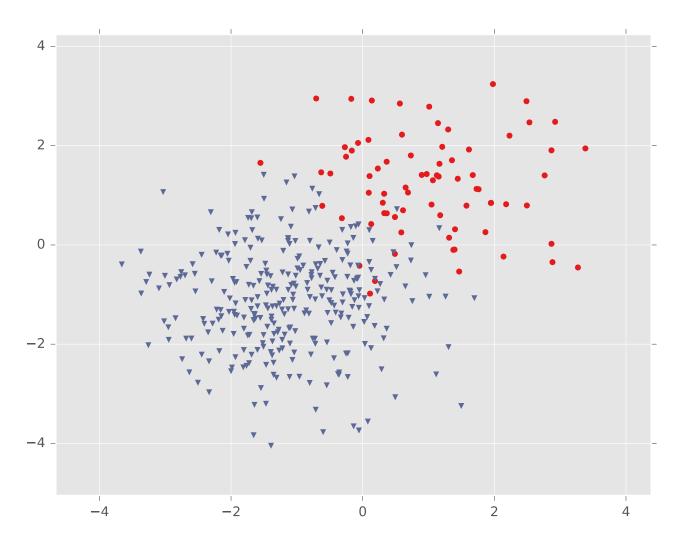
With two input features, $\mathbf{x} = [x_1 \ x_2]^{\mathsf{T}}$, we have two weight parameters and one bias parameter, $\mathbf{w} = [w_1 \ w_2]^{\mathsf{T}}$ and b, that control the slope and vertical offset of the following plane:

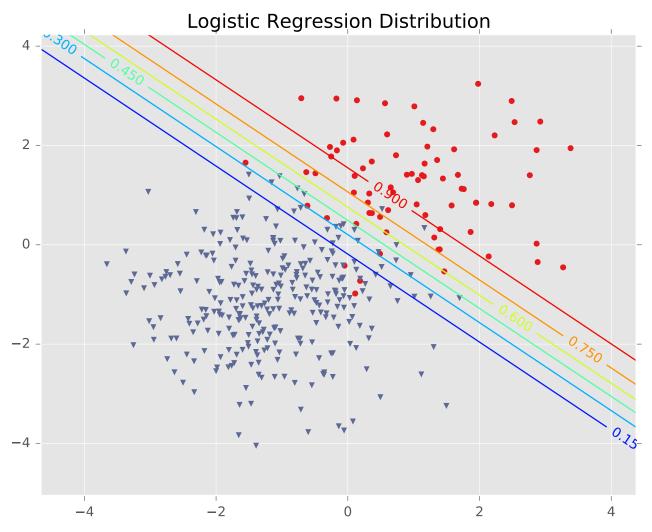
$$z = \mathbf{w}^{\mathsf{T}} \mathbf{x} + b$$

The sigmoid function $\hat{y} = g(z)$ then squashed the plane such that any z values going to $+\infty$ go to 1 and z values going to $-\infty$ go to -1

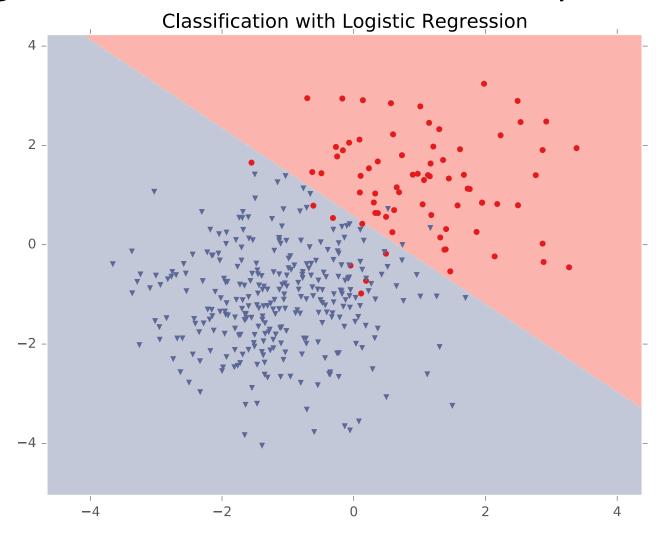






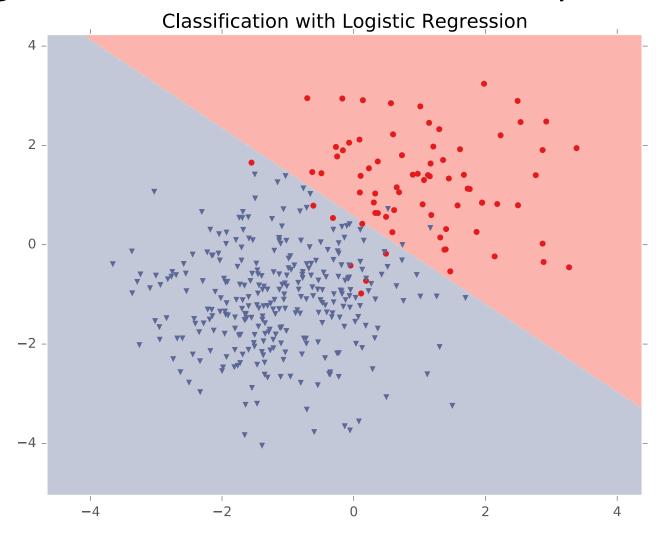


Logistic Regression Decision Boundary



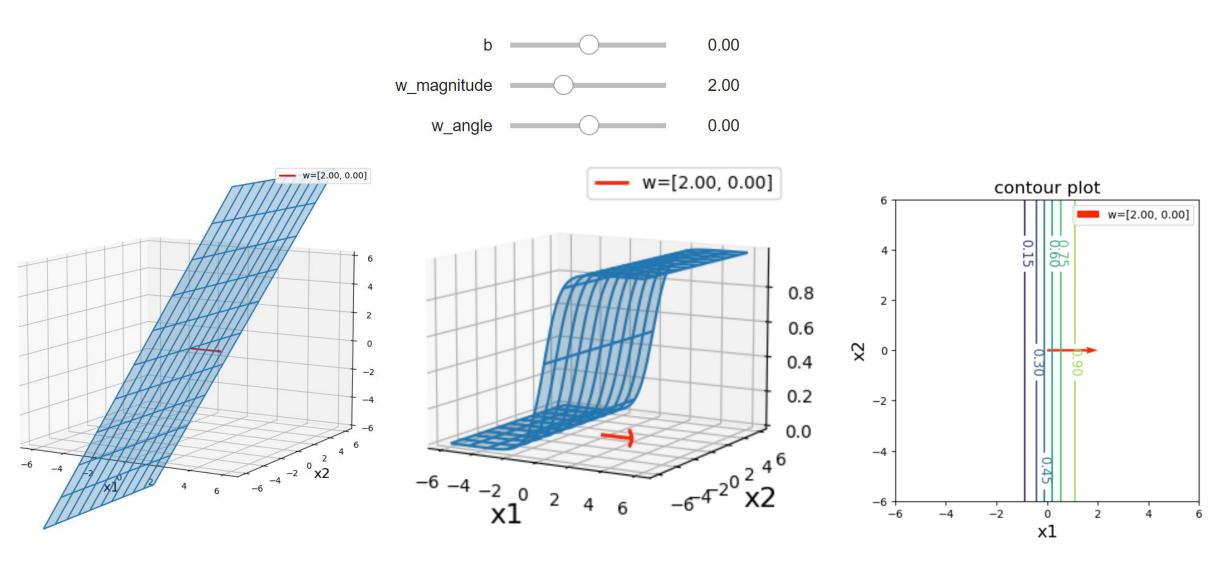
Linear decision boundary

Logistic Regression Decision Boundary



Exercise

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



Linear in Higher Dimensions

point

 $\mathbf{w}^T \mathbf{x} + b = 0$

1. D
$$y = W \times Jb$$

2-D $y = V_1 \times_1 + W_2 \times_2 + b$

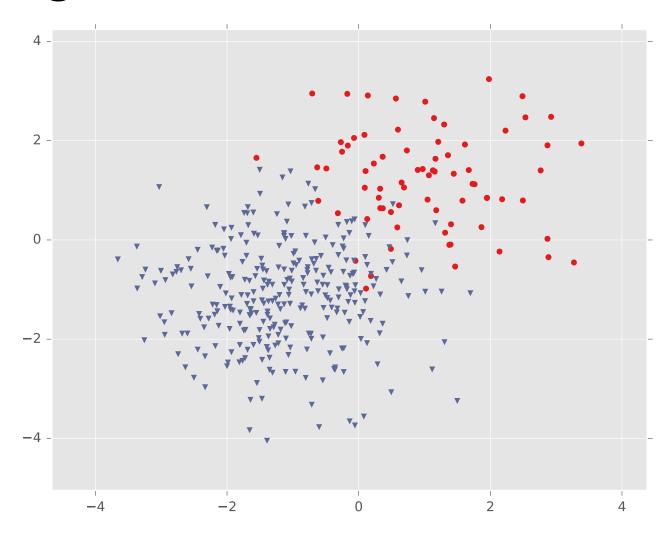
What are these linear shapes called for 1-D, 2-D, 3-D, M-D input?

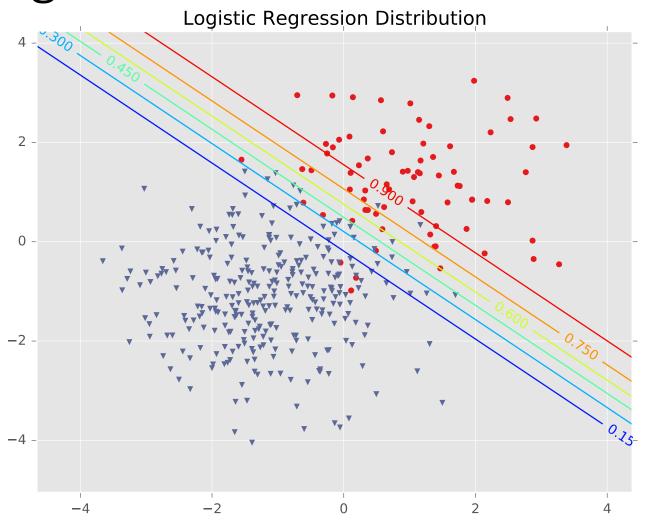
$$x \in \mathbb{R}$$
 $x \in \mathbb{R}^2$ $x \in \mathbb{R}^3$ $x \in \mathbb{R}^M$

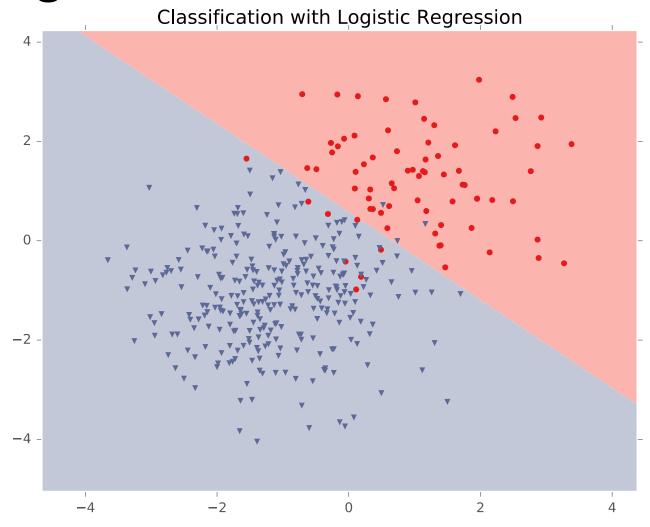
$$y = w^T x + b$$
 line plane hyperplane hyperplane

line

$$w^T x + b \ge 0$$
 halfline halfplane halfspace halfspace





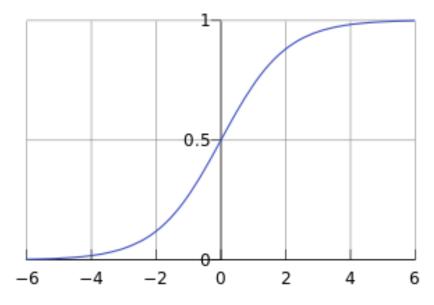


Poll 3

For a point \mathbf{x} on the decision boundary of logistic regression, does $g(\mathbf{w}^T\mathbf{x} + b) = \mathbf{w}^T\mathbf{x} + b$?

- A) Yes
- B) No
- C) I have no idea

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

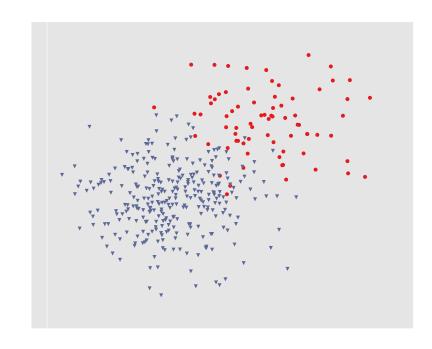


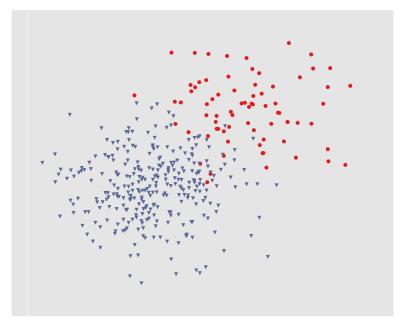
Optimization

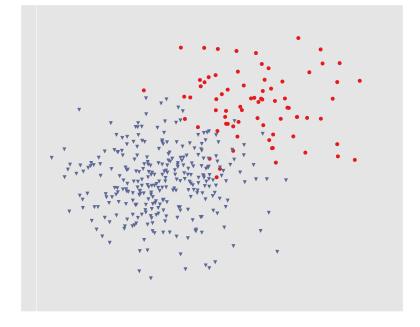
Optimizing a Model for Cancer Diagnosis

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} .

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





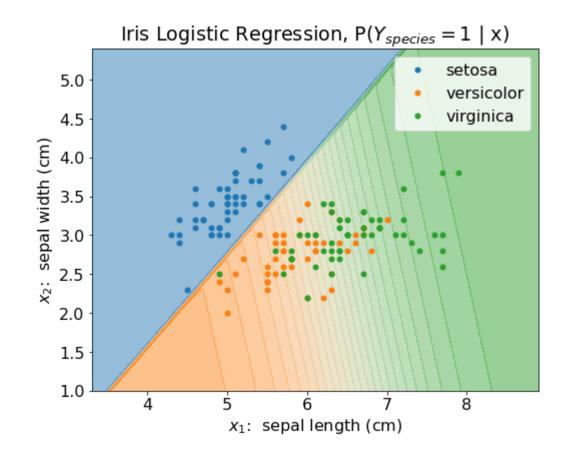


Empirical Risk Minimization

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

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

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



Empirical Risk Minimization

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

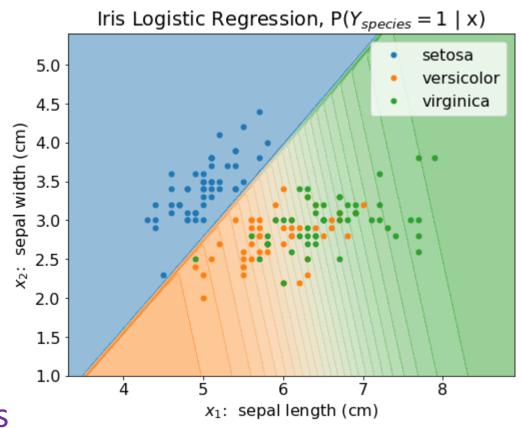
$$h^* = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \ \hat{R}(h)$$

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

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



Binary Logistic Regression

1) Model

2) Objective function

3) Solve for $\widehat{\boldsymbol{\theta}}$

Binary Logistic Regression

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

Objective: Special case for binary logistic regression

$$\hat{y} = g(\boldsymbol{\theta}^T \mathbf{x})$$

$$J(\boldsymbol{\theta}) = -\frac{1}{N} \sum_{i} \sum_{k} y_k^{(i)} \log y_k^{(i)}$$

$$= -\frac{1}{N} \sum_{i} (y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}))$$

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

$$J^{(i)}(\boldsymbol{\theta}) = -[y^{(i)}\log \hat{y}^{(i)} + (1 - y^{(i)})\log(1 - \hat{y}^{(i)})]$$

$$\frac{\partial J^{(i)}}{\partial \boldsymbol{\rho}} = -(y^{(i)} - \hat{y}^{(i)}) \mathbf{x}^{(i)}$$

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

$$J^{(i)}(\boldsymbol{\theta}) = -[y^{(i)}\log \hat{y}^{(i)} + (1 - y^{(i)})\log(1 - \hat{y}^{(i)})]$$

$$\frac{\partial J^{(i)}}{\partial \boldsymbol{\rho}} = -(y^{(i)} - \hat{y}^{(i)}) \mathbf{x}^{(i)}$$

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

$$Z = \partial T \times = \int du \quad uTV$$

$$\hat{y} = g(z)$$

$$= V \quad or \quad VT$$

$$\frac{dg}{dz} = g(z)(1 - g(z))$$

$$J^{(i)}(\boldsymbol{\theta}) = -\left[y^{(i)}\log\hat{y}^{(i)} + (1 - y^{(i)})\log(1 - \hat{y}^{(i)})\right] \mathcal{L}(y^{(i)})$$

$$J = \frac{1}{2} \frac{1}{2}$$

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

$$Z = \partial T \times = \int u V$$

$$\hat{y} = g(z)$$

$$\frac{dg}{dz} = g(z)(1 - g(z))$$

$$= V \text{ or } V^{T}$$

$$J^{(i)}(\boldsymbol{\theta}) = -\left[y^{(i)}\log\hat{y}^{(i)} + (1 - y^{(i)})\log(1 - \hat{y}^{(i)})\right] \mathcal{J}(y^{(i)})$$

$$= -\left[\frac{y}{\gamma} - \frac{y}{\gamma - \hat{y}}\right] \frac{\partial \hat{y}}{\partial \hat{y}}$$

$$= -\left[\frac{y}{\gamma} - \frac{y}{\gamma - \hat{y}}\right] \frac{\partial \hat{y}}{\partial \hat{y}} \frac{\partial z}{\partial \hat{y}}$$

$$= -\left[\frac{y}{\gamma} - \frac{y}{\gamma - \hat{y}}\right] \frac{\partial \hat{y}}{\partial \hat{y}} \frac{\partial z}{\partial \hat{y}}$$

$$\frac{\partial J^{(i)}}{\partial \hat{y}} = -(y^{(i)} - \hat{y}^{(i)}) \mathbf{x}^{(i)}$$

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

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

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

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

No closed form solution 🕾

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

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

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

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

No closed form solution 🕾

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

Good news: The logistic regression optimization function is convex!

Logistic Regression

Convexity

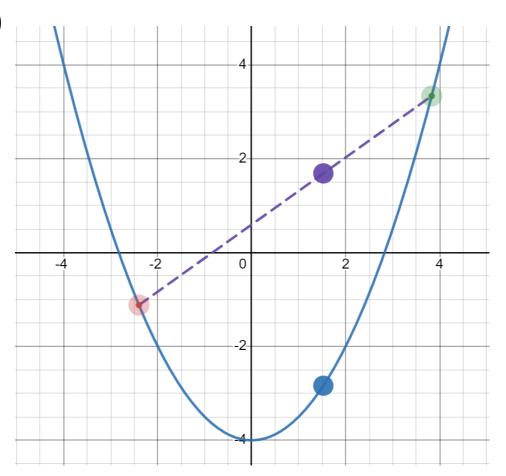
Optimization

Convex function

If f(x) is convex, then:

 $f(\alpha \mathbf{x} + (1 - \alpha)\mathbf{z}) \le \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{z})$ $\forall 0 \le \alpha \le 1$

Demo on Desmos



Optimization

Convex function

If f(x) is convex, then:

•
$$f(\alpha \mathbf{x} + (1 - \alpha)\mathbf{z}) \le \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{z}) \quad \forall \ 0 \le \alpha \le 1$$

Convex optimization

If second derivative is ≥ 0 everywhere then function is convex

If f(x) is convex, then:

■ Every local minimum is also a global minimum ©

Optimization

 $h(\mathbf{x}) = g(\mathbf{w}^{\mathsf{T}}\mathbf{x} + b)$ is definitely not convex

But...what are we optimizing over in logistic regression?

$$J(\theta) = -\frac{1}{N} \sum_{i} \sum_{k} y_{k}^{(i)} \log y_{k}^{(i)}$$

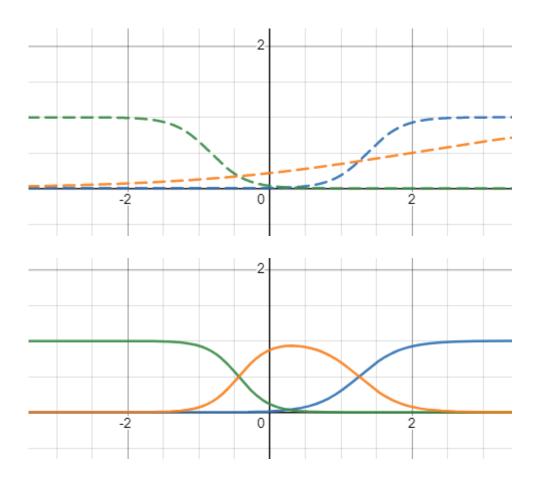
$$= -\frac{1}{N} \sum_{i} (y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}))$$

Multi-class Logistic Regression

Multi-class Logistic Regression

Desmos Demo:

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



Multi-class Logistic Regression

Cross-entropy loss

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

Model

$$\hat{\mathbf{y}} = h(\mathbf{x}) = g_{softmax}(\mathbf{z})$$

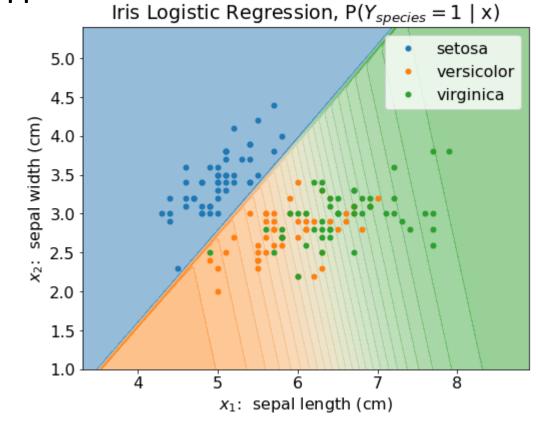
$$z = \Theta x$$

 $z_k = \boldsymbol{\theta}_k \mathbf{x}$

$$\mathbf{x} = \begin{bmatrix} 1 \\ \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}$$

One vector of parameters for each class

$$\boldsymbol{\theta}_k = \begin{bmatrix} b_k \\ w_{k,1} \\ w_{k,2} \end{bmatrix}$$



Stacked into a matrix of $K \times M$ parameters

$$\Theta = \begin{bmatrix} - & \boldsymbol{\theta}_{1}^{\mathsf{T}} - \\ - & \boldsymbol{\theta}_{2}^{\mathsf{T}} - \\ - & \boldsymbol{\theta}_{3}^{\mathsf{T}} - \end{bmatrix} = \begin{bmatrix} b_{1} & w_{1,1} & w_{1,2} \\ b_{2} & w_{2,1} & w_{2,2} \\ b_{3} & w_{3,1} & w_{3,2} \end{bmatrix}$$

Logistic Function

Logistic (sigmoid) function converts value from $(-\infty, \infty) \to (0, 1)$ $g(z) = \frac{1}{1 + e^{-z}} = \frac{e^z}{e^z + 1}$

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

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

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.

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

$$\begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_K \end{bmatrix} \to \begin{bmatrix} e^{z_1} \\ e^{z_2} \\ \vdots \\ e^{z_K} \end{bmatrix} \cdot \frac{1}{\sum_{k=1}^K e^{z_k}} \quad \text{Example} \begin{bmatrix} -1 \\ 4 \\ 1 \\ -2 \\ 3 \end{bmatrix} \to \begin{bmatrix} 0.0047 \\ 0.7008 \\ 0.0349 \\ 0.0017 \\ 0.2578 \end{bmatrix}$$

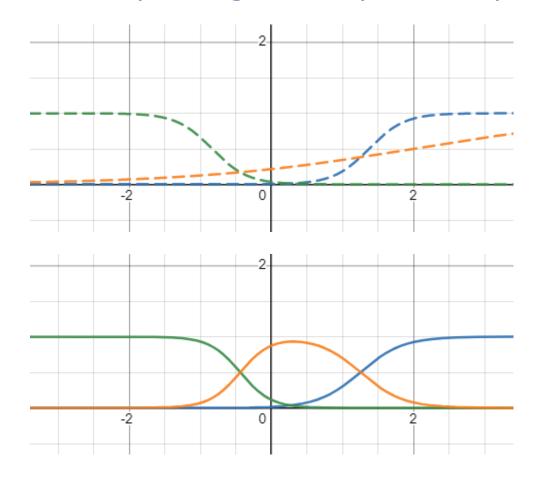
Multiclass Predicted Probability

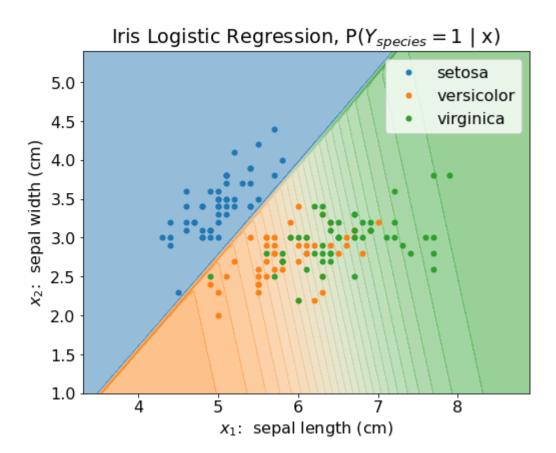
Multiclass logistic regression uses the parameters learned across all K classes to predict the discrete conditional probability distribution of the output Y given a specific input vector \mathbf{x}

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

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}

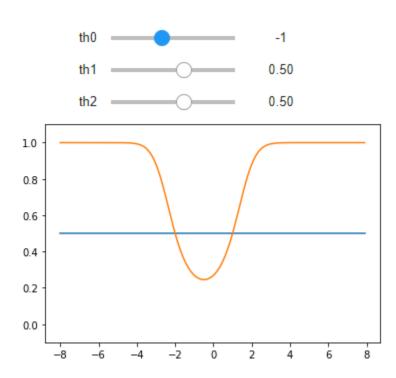




Logistic Regression with Polynomial Features

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

