

# Data X

Data-X : Machine Learning (Classification)

Ikhlaq Sidhu, Kevin Bozhe Li

Ikhlaq Sidhu  
Chief Scientist & Founding Director,  
Sutardja Center for Entrepreneurship & Technology  
IEOR Emerging Area Professor Award, UC Berkeley

# Agenda

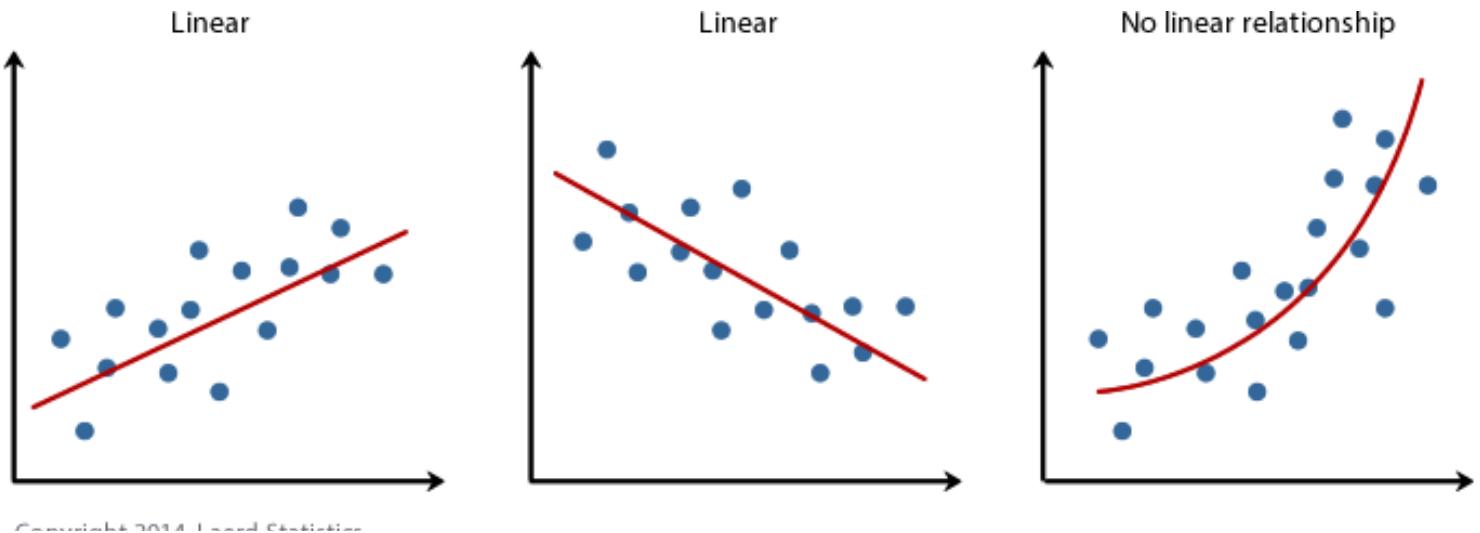
- Discussion of Classification, Logistic Regression, and how it works in SciKit Learn (45 min to 1 hour)
- Break Out Session with Sample Code (40 min)
- Homework for next week: Lighter to make more progress on project
- Project Discussion – Next steps, work with team:
  - User requirements: what you know and don't?
  - Technical needs: what you know and don't
  - Whiteboard Your Architecture and /or Data Model
  - Basic planning and prioritization

Data X



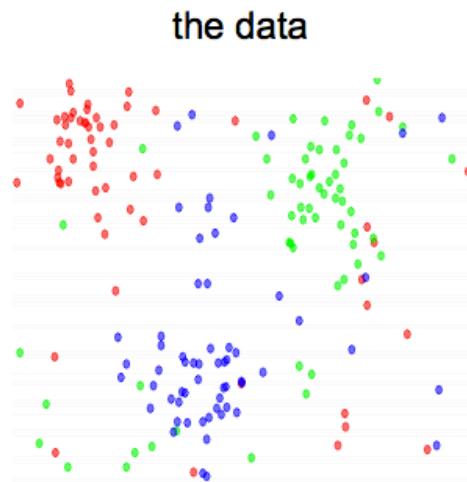
# Prediction vs Classification

Prediction:  
Linear or  
Non-linear

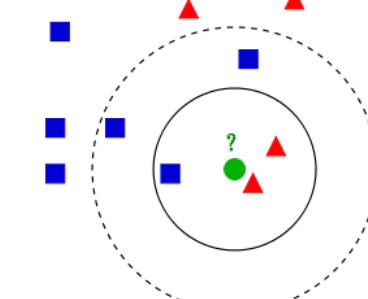


Classification:  
Can also be  
Linear or  
Non-linear

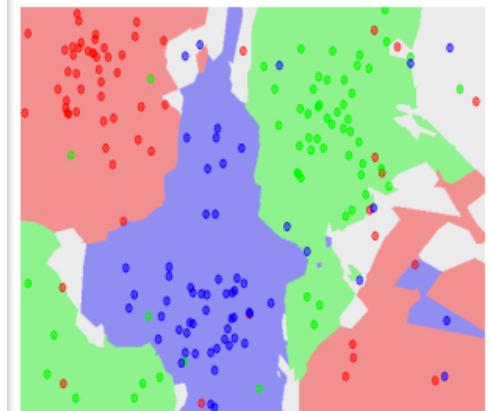
Ex. KNN, K-means  
Logistic, SVM, ..



**KNN Method:** Find the k nearest images and have them vote on the label (i.e. take the mode)



5-NN classifier

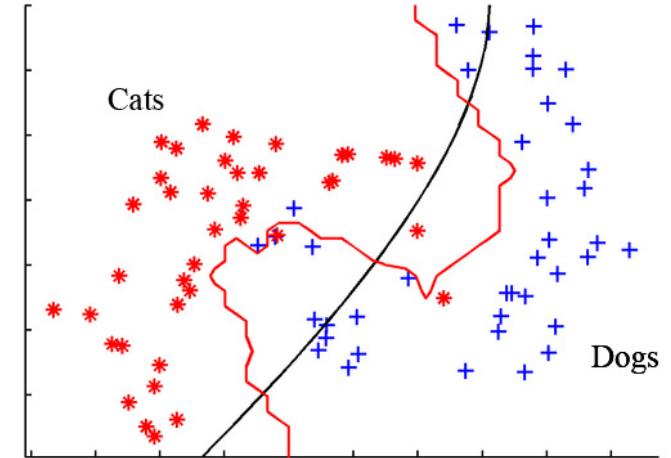


Data X

# Classification

## Examples

- Radar: WWII
- Spam Detection
- Image Classification: Cats VS Dogs
- Image Classification: Recognizing Digits



80322-4129 80206  
40004 14310  
37878 05753  
3502 75316  
35460 44209

Data	X
0	1
1	0
1	1

# Classifiers: A Whirlwind Tour

- **Linear Classifiers**

- The simplest **parametric** classifier
- One must define a **loss function** (a way to measure the error) in order to optimize the weights of a linear classifier
- **Regularization** is your friend!
- Examples: *Support Vector Machine (SVM) & Logistic Classifier*

- **Nonlinear Classifiers**

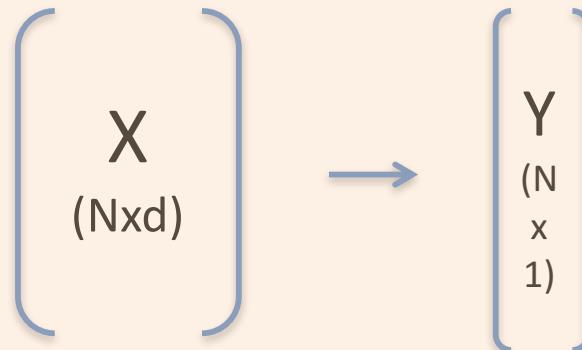
- Introduction to **Neural Networks** (an universal function approximator)



Reminder from Previous Class:  
In the **ML framework**, there is no limit to the predictors or **classifiers** that can be used.

In Sample Data  
which we have

- d features
- N samples



Use for training

Results in  $Y$   
(which we know)  
  
Sometimes measured  
 $Y$  includes error or noise

We don't know:  
  
 $P(X)$ , the distribution of  $X$   
  
Function  $f$ :  
 $f: X \rightarrow Y$

We use this to train:  
 $(x_1, y_1), (x_2, y_2) \dots (x_N, y_N)$

**H = Hypothesis Set**  
All Possible Algorithms  
Candidate Formulas

Algorithm:  
Under test

Final Hypothesis

$g: x \rightarrow y$   
predictor  
or  
classifier

Error measure

- We try different functions  $g$  until
- $g(x)$  is close to  $f(x)$
- For any out-of-sample  $x$ , we can predict  $Y$  or classify it



Data  $X$

# Our Goal: To classify items.



$x_i \xrightarrow{\text{Model: } f(x, W)}$

We have this:

$(X, Y): (x_1, y_1), (x_2, y_2) \dots (x_n, y_n)$

- $x_i$  is a vector (or even matrix) for each data element
- For a picture:  $x_1 = [32 \times 32 \times 3]$ : array of numbers
- $y_i = (-1, +1)$

**Actual Results:**

$$y_i = [y_{i,1}, y_{i,2}, \dots y_{i,k}]$$

$$y_i = [+1, -1, \dots -1]$$

or

$$y_i = [+1, 0, \dots 0]$$

$y(i,1) = 1$  if picture is dog,  
else -1\*

$y(i,2) = 1$  if picture is cat,  
else -1\*

\* Sometimes 0 vs 1  
instead of -1 vs 1

Data X

# Our Goal: To classify items.

We have this: (X,Y)



$x_i \xrightarrow{\text{Model: } f(x,W)}$

**Actual Results:**

$$y_i = [y_{i,1}, y_{i,2}, \dots y_{i,k}]$$

$$y_i = [+1, -1, \dots -1]$$

And sometimes mapped to

$$y_i = [+1, 0, 1..0]$$

## Machine Learning Steps to train a classifier model

1. Choose our **model**:  $f(x_i, W) = \text{our estimate of } Y$

If  $f(x)$  is a Linear model:  $f(i) = f(x_i) = Wx_i = w_1x_{i,1} + w_2x_{i,2} \dots$

$W$  is list of parameters.

Results in a different score for each sample point (picture)  $x_i$

2. Define a **loss function (L)**

3. Optimize across the parameter space ( $W$ ) to minimize the loss function to

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99

Data X

# Example Classification with Logistic Regression

Data: students study for an exam:  
( $x$  = hours studied,  $y$  = pass/not pass)

Hours	0.50	0.75	1.00	1.25	1.50	1.75	1.75	2.00	2.25	2.50	2.75	3.00	3.25	3.50	4.00	4.25	4.50	4.75	5.00	5.50
Pass	0	0	0	0	0	0	1	0	1	0	1	0	1	0	1	1	1	1	1	

In this case output is binary  
(pass or not pass)

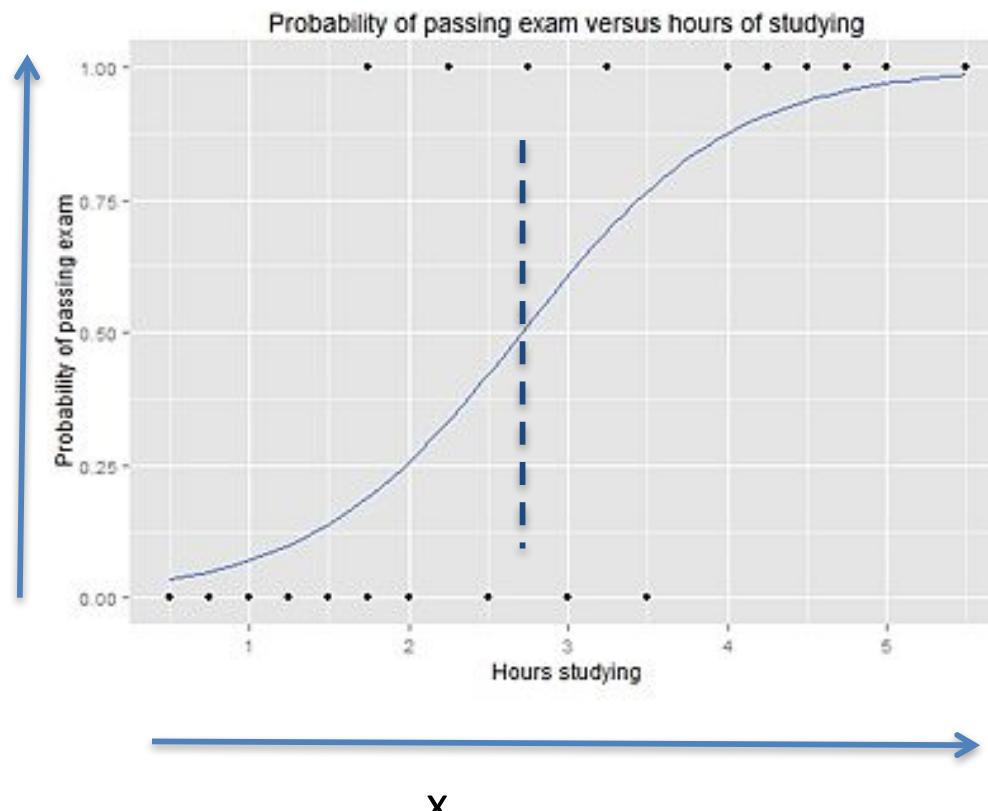
Problem:

Student studies  $x$  hours  
Will the student pass

We use this curve to predict  
the probability that the  
student would pass given  $x$   
hours of study

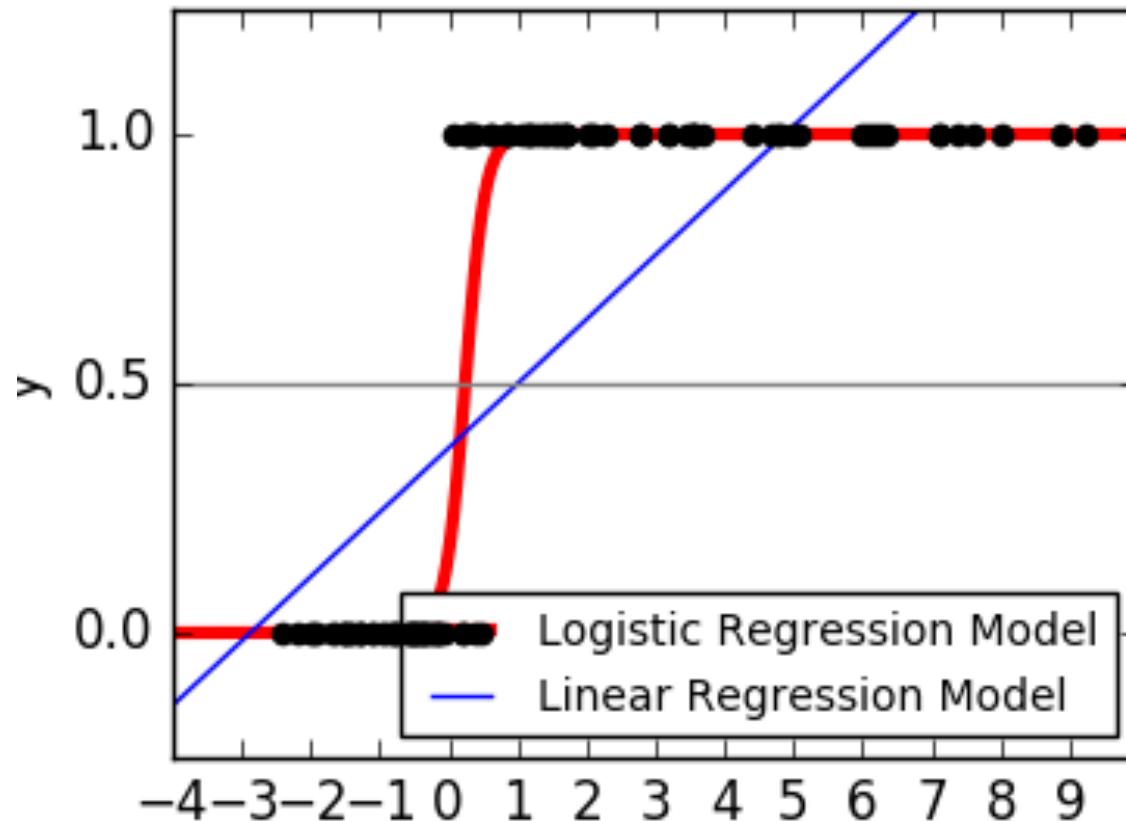
If Prob > 0.5, we classify  
student will pass the exam

$y$  (or t):  
0 = fail  
1 = pass



Data X

Example: Logistic Regression  
We use Logistic Regression because  
a line is not a good estimator for binary results (classification)



Data  $X$

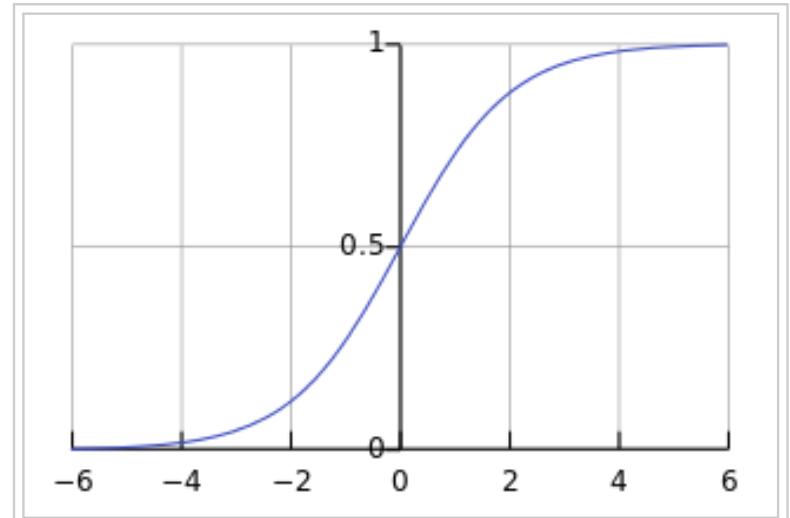
# The logit function is a better fit for binary results

This function is a better fit for binary results:

$$f(t) = \frac{e^t}{e^t + 1} = \frac{1}{1 + e^{-t}}$$

Large  $t \rightarrow 1$   
Small  $t \rightarrow 0$

This function results in (0,1)  
for all real numbers (like a probability)



And if  $t$  has this form:

$$t = w_0 + x_1 w_1$$

$$f(x_i, w_0, w_1) = \frac{1}{1 + e^{-(w_0 + x_1 w_1)}}$$

- Think of  $f(x, w)$  as probability  $y = 1$  given any  $x$
- Prob ( $y=1$ ) =  $\frac{1}{2}$  when  $e^{-(w_0 + x_1 w_1)} = 1$ , ie  $w_0 + x_1 w_1 = 0$
- Choose  $w_0, w_1, \dots, w_d$  to get best fit
- Still need a loss function and then solve for best  $W$



# A Loss Function Commonly Used for Logistic Regression

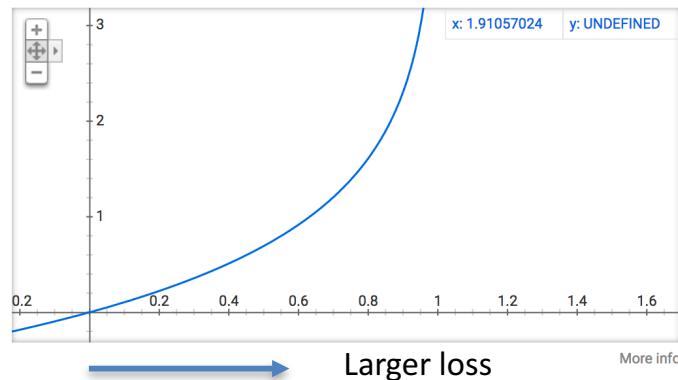
Cross Entropy =  $-t \ln(f(\vec{x})) - (1 - t) \ln(1 - f(\vec{x}))$

If  $y$  is -1 or +1, then we can map to  $t$ : (like actual value  $y$ )

$t = (1+y)/2$ , so  $t$  ranges from 0 to 1

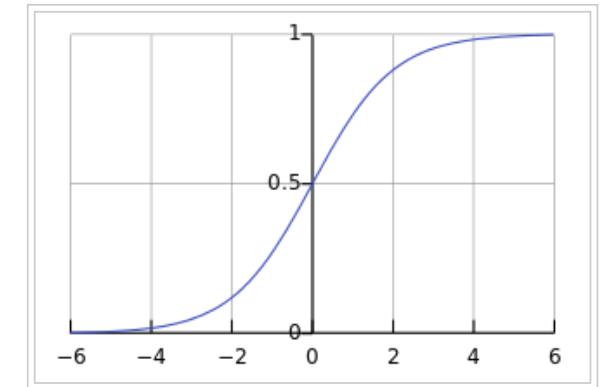
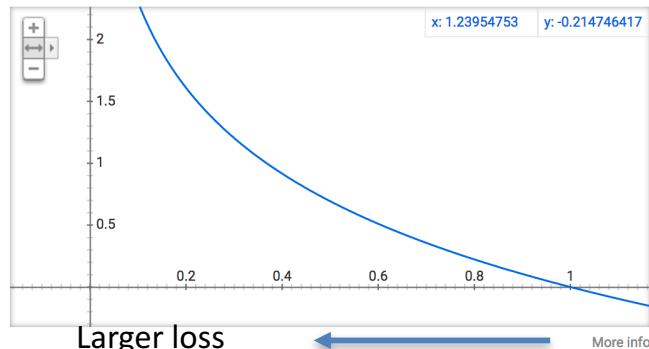
$-\ln(1-f(x))$  when  $t = 0$

if  $t = 0$ : loss =  $-\ln(1-f(x))$   
 if  $f(x)$  is 1  $\rightarrow$  loss =  $e$ ,  
 if  $f(x)$  is 0  $\rightarrow$  loss = 0



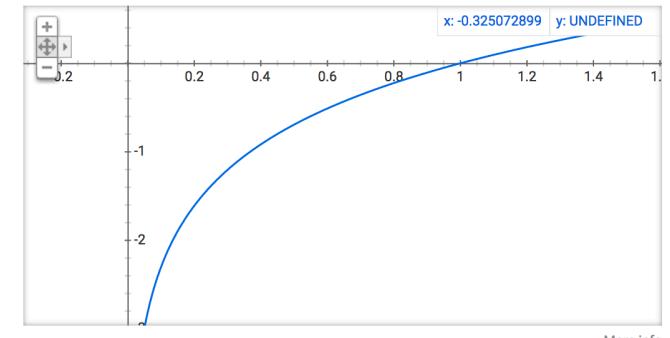
$-\ln(f(x))$  when  $t = 1$

if  $t = 1$ : loss =  $-\ln(f(x))$   
 if  $f(x)$  is near 0  $\rightarrow$  loss = large  
 if  $f(x)$  is near 1  $\rightarrow$  loss = approx. = 0



$$f(x_i, w_0, w_1) = \frac{1}{1 + e^{-(w_0 + x_1 w_1)}}$$

Graph for  $\ln(x)$



Data X

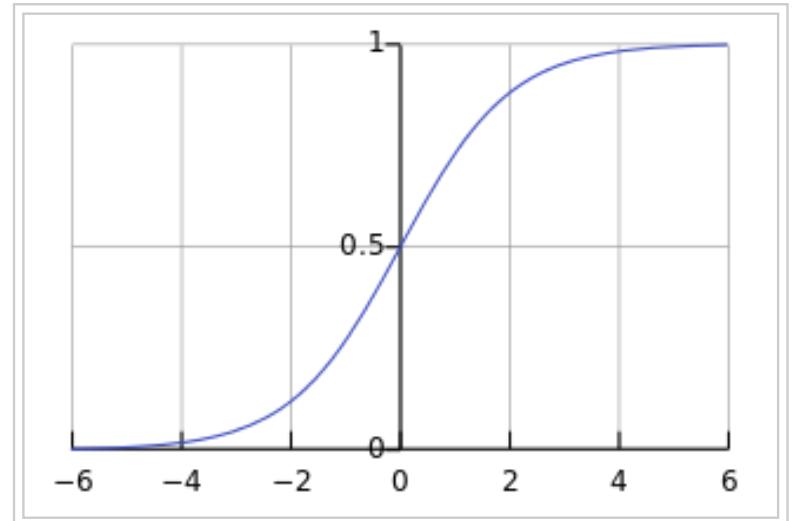
# The logit function can be extended for multiple features (dimensions)

This function is a better fit for binary results:

$$f(t) = \frac{e^t}{e^t + 1} = \frac{1}{1 + e^{-t}}$$

Large  $t \rightarrow 1$   
Small  $t \rightarrow 0$

This function results in (0,1)  
for all real numbers (like a probability)



And if  $t$  has this form:

$$t = w_0 + x_1 w_1$$

or

$$t = Wx_i \text{ (in matrix form)}$$

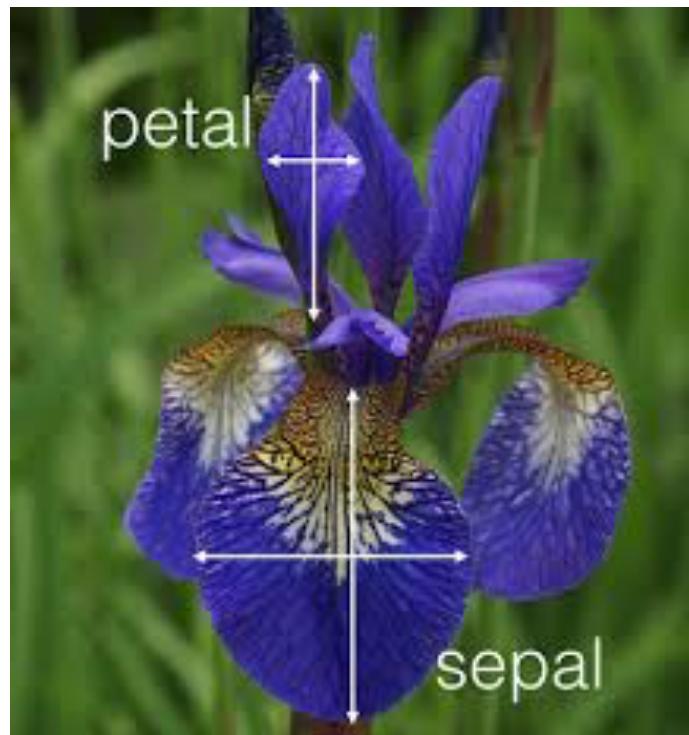
$$= w_1 x_{i,1} + w_2 x_{i,2} \dots$$

- Easily extends to multiple features ( $x_1, x_2, \dots$ )
- And multiple parameter weights

$$f(x_i, W) = \frac{1}{1 + e^{-(w_0 + x_{i,1} w_1 + x_{i,2} w_2 \dots)}}$$



# Example Code Sample with Logistic Regression Classifier



Input data

X: Attribute Information:

- sepal length in cm
- sepal width in cm
- petal length in cm
- petal width in cm

Y:

- 0 = 'setosa',
- 1 = 'versicolor',
- 2 = 'virginica'

```
print type (X)
print X[0:5]
```

```
<type 'numpy.ndarray'>
[[ 5.1  3.5  1.4  0.2]
 [ 4.9  3.  1.4  0.2]
 [ 4.7  3.2  1.3  0.2]
 [ 4.6  3.1  1.5  0.2]
 [ 5.  3.6  1.4  0.2]]
```

```
print Y[0:5]
[0 0 0 0 0]
```

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99

Data X

# Example Code Sample with Logistic Regression Classifier

1 →

```
import numpy as np  
from sklearn import linear_model, datasets  
  
X = iris.data[:, 1:3] # only the first two features.  
Y = iris.target
```

2 →

```
# https://en.wikipedia.org/wiki/Logistic_regression  
logreg = linear_model.LogisticRegression(C=1e5)
```

3 →

```
# we create an instance of Neighbours Classifier and  
fit the data.  
logreg.fit(X, Y)
```

4 →

```
# predict a category for every row in X  
Z = logreg.predict(X)
```

\* Z[2] will be the predicted number for row X[2]

Class  
**sklearn.linear\_model.**  
**LogisticRegression**

(penalty='l2',  
dual=False,  
tol=0.0001,  
**C=1.0**,  
fit\_intercept=True,  
intercept\_scaling=1,  
class\_weight=None,  
random\_state=None,  
solver='liblinear', max\_iter=100,  
multi\_class='ovr', verbose=0,  
warm\_start=False,  
n\_jobs=1)

[http://scikit-learn.org/stable/modules/generated/sklearn.linear\\_model.LogisticRegression.html](http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html)

Data X

```
fit(X, y, sample_weight=None)
```

[source]

Fit the model according to the given training data.

**Parameters:** **X** : {array-like, sparse matrix}, shape (n\_samples, n\_features)

Training vector, where n\_samples is the number of samples and n\_features is the number of features.

**y** : array-like, shape (n\_samples,)

Target vector relative to X.

**sample\_weight** : array-like, shape (n\_samples,) optional

Array of weights that are assigned to individual samples. If not provided, then each sample is given unit weight.

*New in version 0.17: sample\_weight support to LogisticRegression.*

**Returns:** **self** : object

Returns self.

```
predict(X)
```

[source]

Predict class labels for samples in X.

**Parameters:** **X** : {array-like, sparse matrix}, shape = [n\_samples, n\_features]

Samples.

**Returns:** **C** : array, shape = [n\_samples]

Predicted class label per sample.



## Fit and predict from ScikitLearn

# Code Samples with SciKit Learn

```
# Plot the decision boundary. For that, we will assign a color to each
# point in the mesh [x_min, x_max]x[y_min, y_max].
x_min, x_max = X[:, 0].min() - .5, X[:, 0].max() + .5
y_min, y_max = X[:, 1].min() - .5, X[:, 1].max() + .5
xx, yy = np.meshgrid(np.arange(x_min, x_max, h), np.arange(y_min, y_max, h))
Z = logreg.predict(np.c_[xx.ravel(), yy.ravel()])
# numpy.ravel: Return a contiguous flattened array.
```

xx shape is (171, 231)

yy shape is (171, 231)

np.c returns shape (39501, 2)

[[ 3.8 1.5 ]

[ 3.82 1.5 ]

[ 3.84 1.5 ] ...]

Z shape is (39501,)

xx-> X[:,0]	1.5	yy -> X[:,1]	4.9
3.8	3.8, 1.5	3.8, 1.7	3.8, 1.9 ....
4.0	4.0, 1.5	4.0, 1.7	4.0, 1.9....
4.2	4.2, 1.5	4.2, 1.7	4.2, 1.9....
...			
8.4			

Data X

xx is a matrix of all the first values	1.5	4.9
3.8	3.8	3.8 ....
4.0	4.0	4.0 ....
4.2	4.2	4.2....
...		
yy is matrix of only the second values	1.5	4.9
xx-> X[:,0]	1.5	yy -> X[:,1]
3.8	1.5	1.7
4.0	1.5	1.7
4.2	1.5	1.7
...		
8.4		

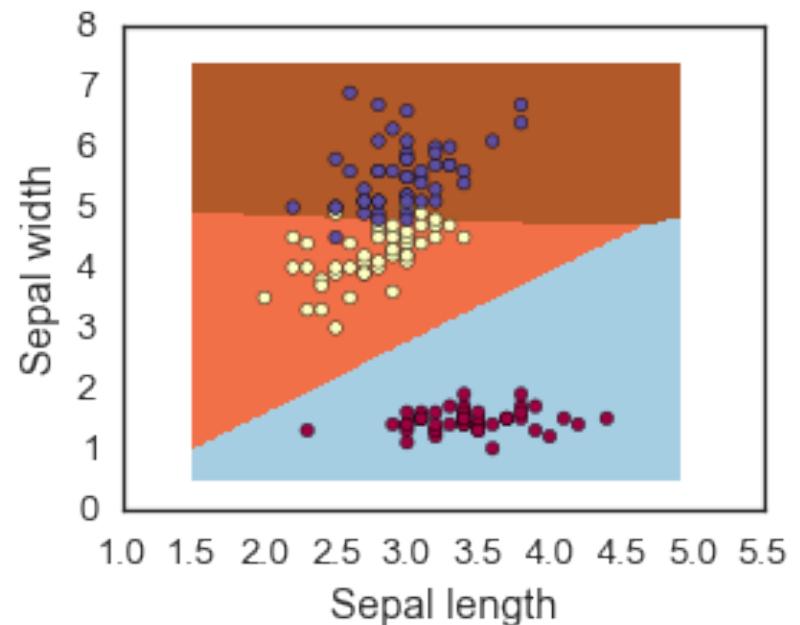
# Plotting the Results

```
# Put the result into a color plot
Z = Z.reshape(xx.shape)
plt.figure(1, figsize=(4, 3))
plt.pcolormesh(xx, yy, Z, cmap=plt.cm.Paired)

# Plot also the training points
plt.scatter(X[:, 0], X[:, 1], c=Y, edgecolors='k',
cmap=get_cmap("Spectral"))
plt.xlabel('Sepal length')
plt.ylabel('Sepal width')

#plt.xlim(xx.min(), xx.max())
#plt.ylim(yy.min(), yy.max())
#plt.xticks(())
#plt.yticks(())

plt.show()
```



# Loss Functions



$x_i = [32 \times 32 \times 3]$ :  
array of numbers

1. If  $f(x)$  is our selected model

If  $f(x)$  is a Linear:

$$\text{Score}(i) = f(x_i, W) = Wx_i$$

Actual Results:

$$y_i = [y_{i,1}, y_{i,2}, \dots y_{i,k}]$$

$$y_i = [+1, -1, \dots -1]$$

Score

Score  $f(x_i, W)$  typically ranges from (-1,+1)

2. Loss Function = price of inaccuracy of score:

- example:  $L_i(f(x_i), W, y) = (1 - y f(\vec{x}))^2$  (square loss)

- See that  $y f(x)$  approx =  $\begin{cases} +1 & \text{for matches: } y=1, f(x)=1 \text{ or } y=-1, f(x) = -1 \\ -1 & \text{for does not match: } y=1, f(x)=-1, y=-1, f(x)=1 \end{cases}$

$L_i(f(x_i, W), y_i)$  = must be selected (there are many options)

Loss  $L = E[L_i]$  = average all  $i$  sample data during training or testing



Data X

# Other Famous Loss Functions

- Logistic Loss  $= \frac{1}{\ln 2} \ln(1 + e^{-yf(\vec{x})})$   
 $y f(x) = 1$  for match: proportional to  $\exp(1+e^{-1})$   
 $y f(x) = -1$  for no match, proportional to  $\exp(1+e^1)$
- Cross Entropy  $= -t \ln(f(\vec{x})) - (1 - t) \ln(1 - f(\vec{x}))$   
 $t = (1+y)/2$ , so  $t$  ranges from 0 to 1

Note: Used for  
Logistic Regression

if  $t = 0$ : loss =  $-\ln(1-f(x)) = e$ , when  $f(x)$  is 1

If  $t = 1$ : loss =  $-\ln(f(x))$  = large when score  $f(x)$  is near 0  
approx. = 0, when  $f(x) = 1$



# After we have a loss function, we need to choose W (parameters) to min the loss

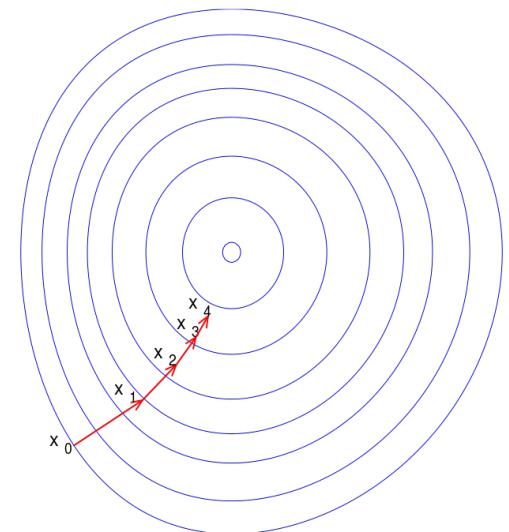
$L(X, W, Y)$  = Loss is a function of  $f(x, W)$  and  $Y$

ML algorithms try a random  $W$

Expected Loss is calculated for small changes in  $W$

$$W_{i+1} = W_i - [(\text{step size gamma}) \times \text{-gradient of } L(W)]$$

Gradient Descent



Gradient descent is a first-order iterative optimization algorithm. To find a local minimum of a function using gradient descent, one takes steps proportional to the negative of the gradient (or of the approximate gradient) of the function at the current point.

Data X

End of Section

Data X

# Regularization

**Why:** To avoid over-fitting

**How:** You penalize your loss function by adding a multiple of an L1 (LASSO) or an L2 (Ridge) norm of your weights vector  $w$

Your new loss function =  $L(X,Y) + \lambda N(w)$

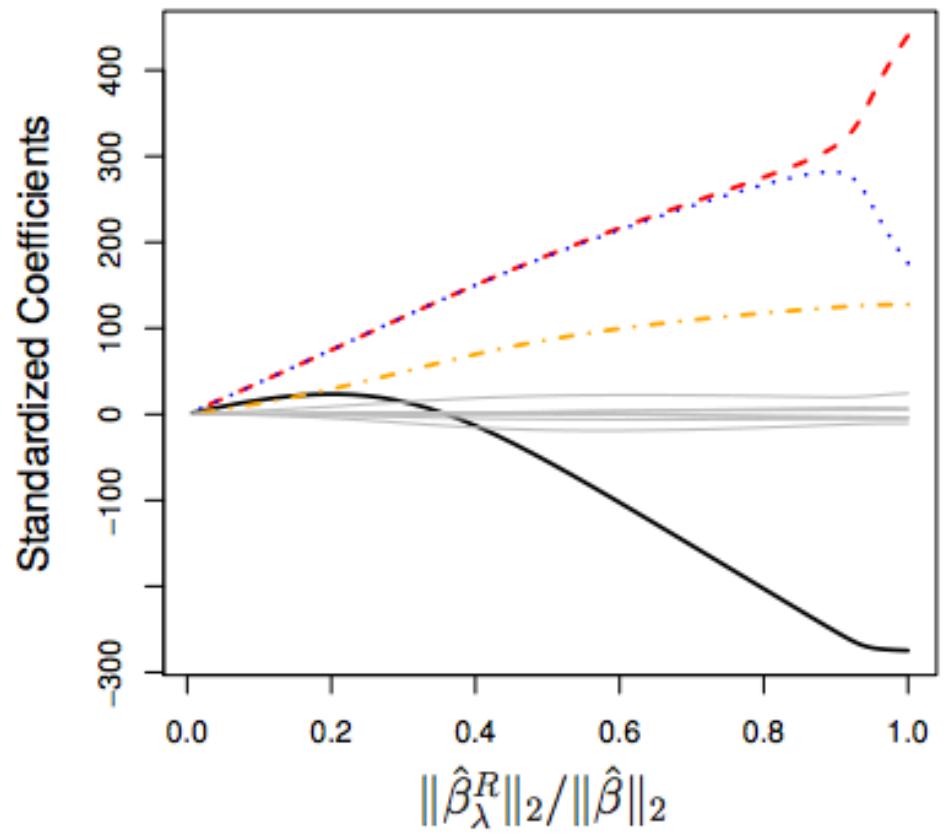
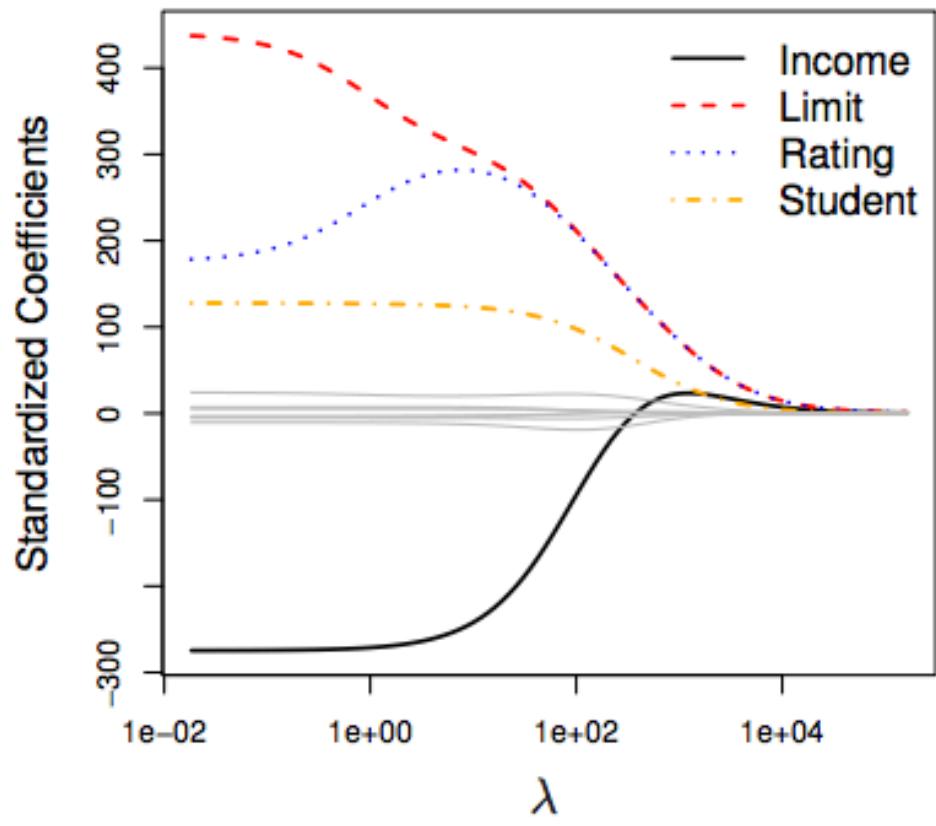
**Tuning the regularization term  $\lambda$ :** Cross-validation:

- divide your training data,
- train your model for a fixed value of  $\lambda$  and test it on the remaining subsets
- repeat this procedure while varying  $\lambda$ .

Then you select the best  $\lambda$  that minimizes your loss function.

Data X

# Shrinkage Methods II: An example



# References

- The material presented in this lecture references lecture material draws on the materials the following courses:
- UC Berkeley – CS 294-129 (Designing, Visualizing, and Understanding Deep Neural Networks):  
<https://bcourses.berkeley.edu/courses/1453965/pages/cs294-129-designing-visualizing-and-understanding-deep-neural-networks>
- Stanford – CS231n (Convolutional Neural Networks for Visual Recognition):  
<http://cs231n.stanford.edu/>

