

# Lecture 1: Machine Learning Basics

Siddharth Garg  
sg175@nyu.edu

# This Course...



Social network  
deanonymization

Spam filtering

Growing use of ML  
techniques in cyber-  
security application

Biometrics

Browser  
fingerprinting

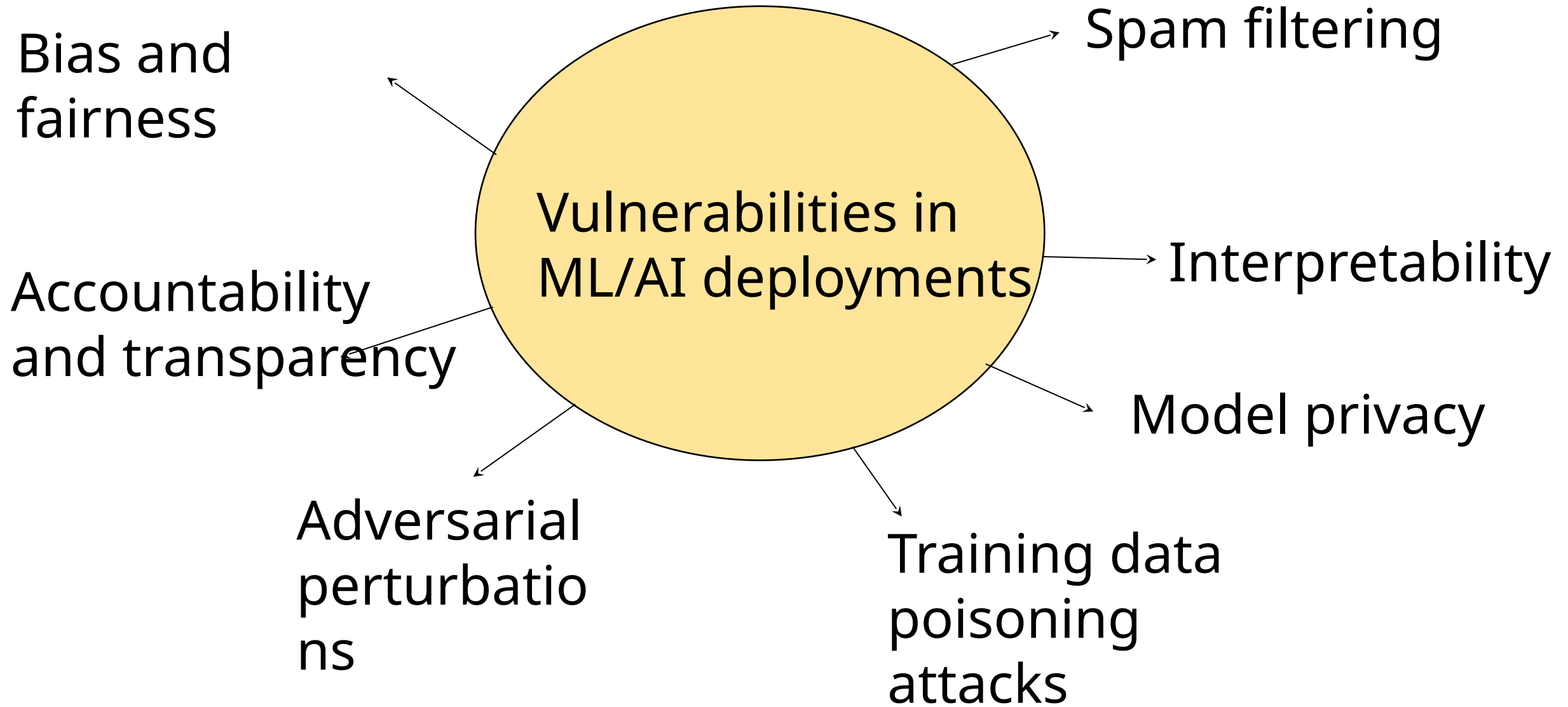
Malware  
detection



Automated  
Evasion

Network  
intrusion  
detection

# This Course...



# What is Machine Learning?

- Ability for machines to learn without being *explicitly* programmed

"A computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$  if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ ." --- Mitchell, T. (1997). Machine Learning. McGraw Hill. p. 2.

- Why not use user knowledge, experience or expertise?
  - Are humans always able to explain their expertise?
  - Can machines *outperform* humans?
- What kinds of experiences ( $E$ ), tasks ( $T$ ) and performance measures ( $P$ )?

# Example: MNIST Digit Recognition

## Task

(T): • Given gray-scale images  $x \in [0,255]^{28 \times 28}$  and  $y \in [0,9]$  find a function

$$f: x \rightarrow y$$



<https://www.npmjs.com/package/mnist>

## Experience

(E): A "training dataset" a set of correctly labeled images

## Performance

(P): Accuracy on a "test dataset"

**"Supervised Learning (Classification)"**

# Example: Spam Classification

## Task

(T): • Emails  $x \in$  all possible emails and  $y \in \{spam, non\_spam\}$  find

$$f: x \rightarrow y$$

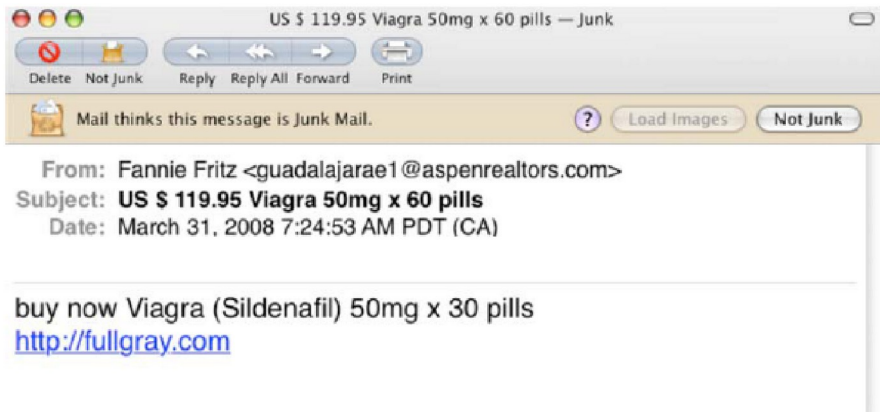
## Experience

(E): A "training dataset" a emails marked as "spam" or "non\_spam"

## Performance

(P): Spam detection accuracy

"Supervised Learning (Classification)"



↓  
SPAM

# Some Challenges

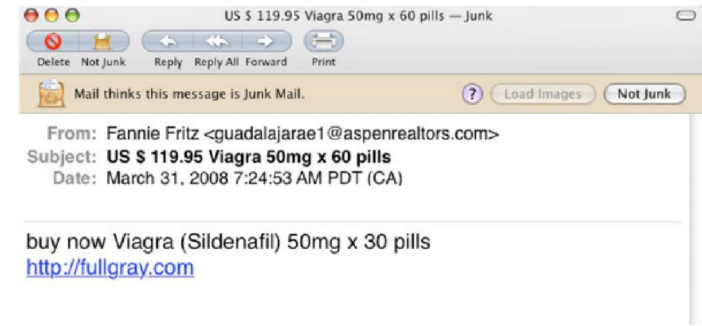
## Representing Data (or Feature Extraction)

How to represent  $x \in$  all possible emails *mathematically*

- One example is “**bag of words**” representation: # times each word in the dictionary occurs
  - What do you lose?
  - What do you gain?
  - How can we compress this representation further?

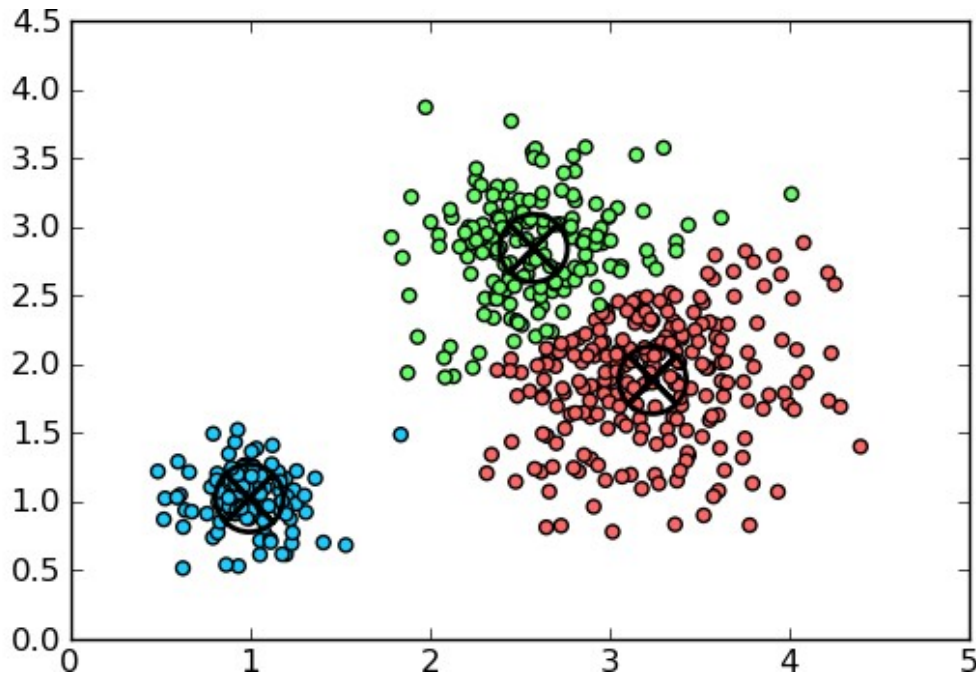
## What kind of classifier?

- What does the function  $f$  look like?
- And how do we learn its parameters?



# Example: Clustering

**Task (T):** “Cluster” a set of documents into  $k$  groups such that “similar” documents appear in the same group



## Experience

**(E):** A “training dataset” of documents without “labels”

## Performance

**(P):** Average distance to cluster center

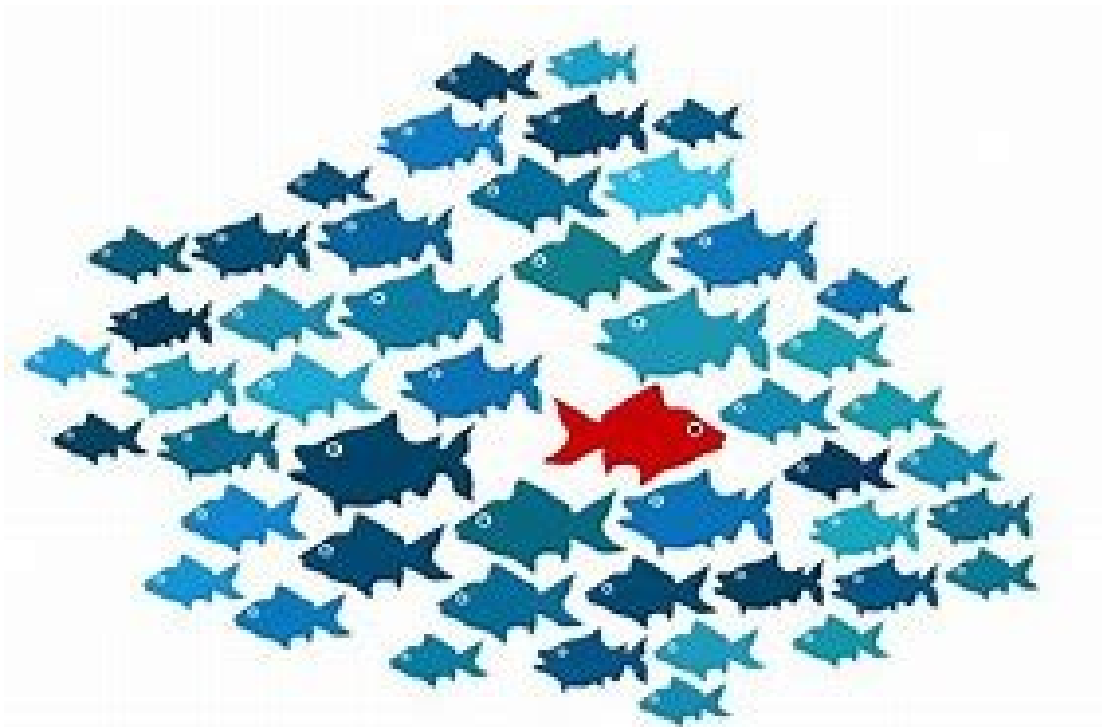
“Unsupervised Learning”



# Example: Anomaly Detection

## Task

(T): • Which of these is like the others?



## Experience

(E): Unlabeled samples

## Performance

(P): Anomaly detection accuracy

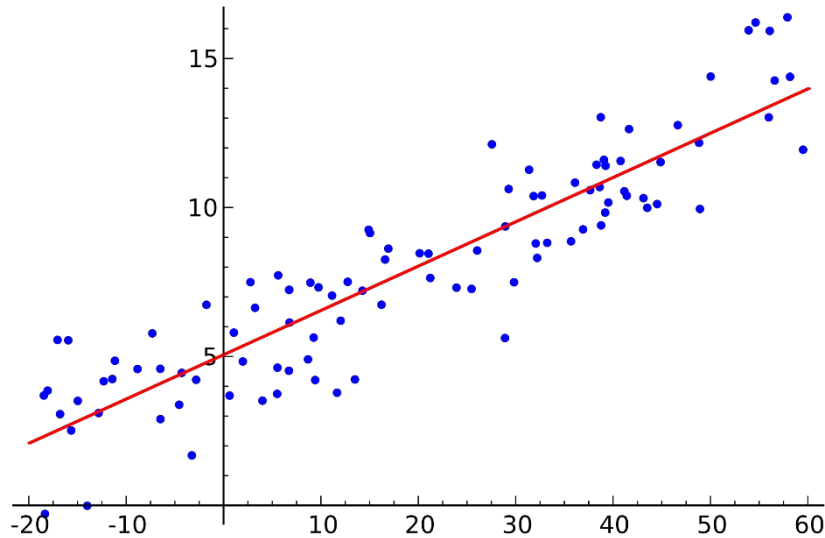
*“Unsupervised  
Learning”*

# Regression

## Task

**(T):** • Given  $x \in \mathbb{R}$  and  $y \in \mathbb{R}$  find a *linear* function

$$f: x \rightarrow y$$



[S. Rangan, EL-GY-9123  
Lec 2]

## Experience

**(E):** Training data: Points  $(x_i, y_i), i \in [1, N]$

## Performance

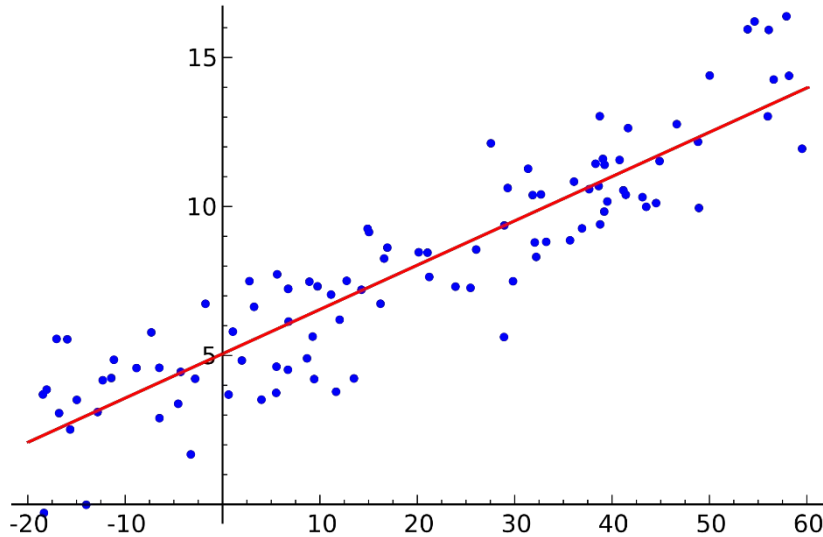
**(P):** *Least squares fit*: minimize mean square error between prediction and ground-truth

***"Supervised Learning  
(Regression)"***

# Linear Least Squares Regression

$$y = f(x) = \beta_1 x + \beta_0$$

- How do we find the values  $(\beta_1, \beta_0)$ ?



$$\min_{\beta_1, \beta_0} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

$$\hat{y}_i = \beta_1 x_i + \beta_0 \quad \forall i \in [1, N]$$

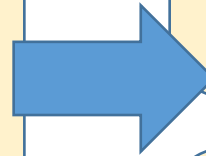
# Linear Least Squares Regression

$$y = f(x) = \beta_1 x + \beta_0$$

- How do we find the values  $(\beta_1, \beta_0)$ ?

$$\min_{\beta_1, \beta_0} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

$$\hat{y}_i = \beta_1 x_i + \beta_0 \quad \forall i \in [1, N]$$



$$\min_{\beta_1, \beta_0} \sum_{i=1}^N (y_i - \beta_1 x_i - \beta_0)^2$$

$$g(\beta_1, \beta_0)$$



$$\frac{\partial g}{\partial \beta_1} = 0 \qquad \frac{\partial g}{\partial \beta_0} = 0$$

# Linear Least Squares Regression

$$y = f(x) = \beta_1 x + \beta_0$$

- How do we find the values  $(\beta_1, \beta_0)$ ?

Residual Sum Squares

$g(\beta_1, \beta_0)$

(RSS)

$$\min_{\beta_1, \beta_0} \sum_{i=1}^N (y_i - \beta_1 x_i - \beta_0)^2$$

$$\frac{\partial g}{\partial \beta_1} = 0$$

$$\frac{\partial g}{\partial \beta_0} = 0$$

$$\frac{\partial g}{\partial \beta_0} = \sum_{i=1}^N -2(y_i - \beta_1 x_i - \beta_0) = 0$$

Sample mean

$$\beta_0 = \frac{\sum_{i=1}^N (y_i - \beta_1 x_i)}{N} = \bar{y} - \beta_1 \bar{x}$$

Are you surprised?

# Linear Least Squares Regression

- How do we find the values  $(\beta_1, \beta_0)$ ?

$$\min_{\beta_1, \beta_0} \sum_{i=1}^N (y_i - \beta_1 x_i - \beta_0)^2 \quad g(\beta_1, \beta_0)$$

$$\frac{\partial g}{\partial \beta_1} = 0$$

$$\frac{\partial g}{\partial \beta_0} = 0$$

$$\frac{\partial g}{\partial \beta_1} = \sum_{i=1}^N -2x_i(y_i - \beta_1 x_i - \beta_0) = 0$$

$$\sum_{i=1}^N x_i((y_i - \bar{y}) - \beta_1(x_i - \bar{x})) = 0$$

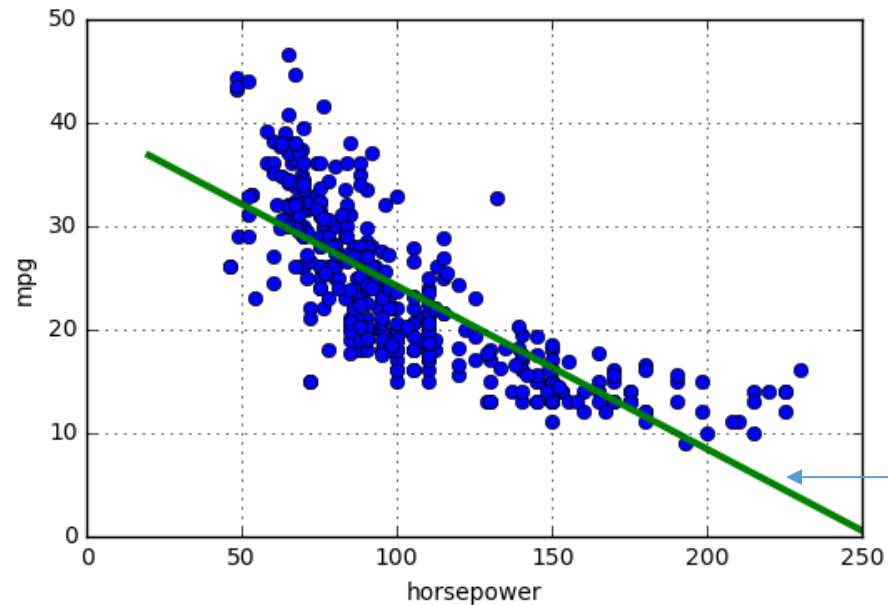
$$\beta_1 = \frac{\overline{xy} - \bar{x}\bar{y}}{\overline{x^2} - \bar{x}^2}$$

Sample covariance

Sample variance

# Auto Example

- Python code



```
xm = np.mean(x)
ym = np.mean(y)
syy = np.mean((y-ym)**2)
syx = np.mean((y-ym)*(x-xm))
sxx = np.mean((x-xm)**2)
beta1 = syx/sxx
beta0 = ym - beta1*xm
```

beta0= 39.94, beta1= -0.16

Regression line:

$$\text{mpg} = \beta_0 + \beta_1 \text{ horsepower}$$

# Linear Least Squares (Multivariate)

- Now consider input  $x \in \mathbb{R}^M$  and output  $y \in \mathbb{R}$  the goal is to learn

$$y = f(x) = \beta_M x_M + \dots + \beta_1 x_1 + \beta_0$$

- Given training dataset  $X \in \mathbb{R}^{N \times M}$  and  $Y \in \mathbb{R}^N$

Training  
sample

$$\begin{matrix} \rightarrow \\ \begin{pmatrix} 1 & x_{01} & x_{02} & \dots & x_{0M} \\ 1 & x_{11} & x_{12} & \dots & x_{1M} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N1} & x_{N2} & \dots & x_{NM} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_M \end{pmatrix} = \begin{pmatrix} y_0 \\ y_1 \\ \vdots \\ y_N \end{pmatrix} \end{matrix}$$

$$\hat{Y} = X\beta$$

Note: for simplicity we will assume that  $X$  includes a column of 1s



# Linear Least Squares (Multivariate)

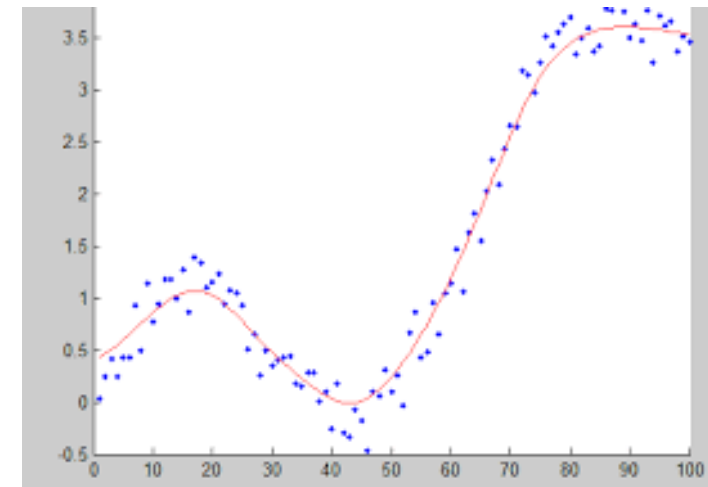
$$RSS = \sum (y - \hat{y})^2 = (Y - \hat{Y})^T \times (Y - \hat{Y}) = (Y - X\beta)^T \times (Y - X\beta)$$

**Objective:**  $\min_{\beta} (Y - X\beta)^T \times (Y - X\beta)$

**Solution:**  $\beta^* = (X^T X)^{-1} X^T Y$

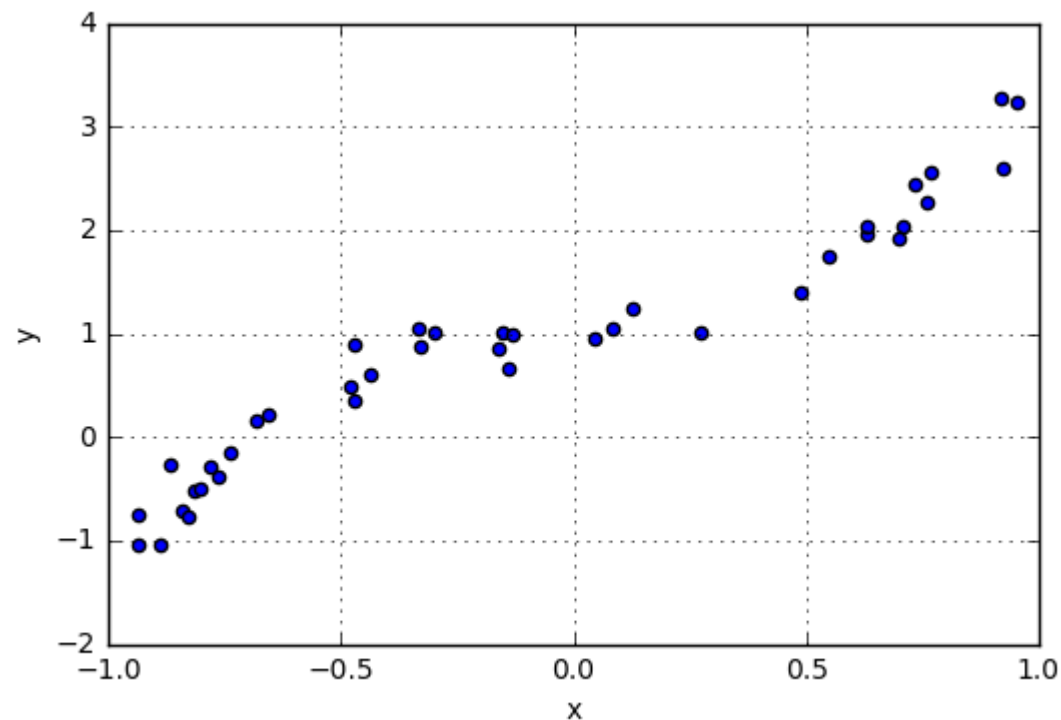
# Polynomial Fitting

- Last lecture: polynomial regression
- Given data  $(x_i, y_i), i = 1, \dots, N$
- Learn a polynomial relationship:
$$y = \beta_0 + \beta_1 x + \dots + \beta_d x^d + \epsilon$$
  - $d$  = degree of polynomial. Called **model order**
  - $\boldsymbol{\beta} = (\beta_0, \dots, \beta_d)$  = coefficient vector
- Given  $d$ , can find  $\boldsymbol{\beta}$  via least squares
- How do we select  $d$  from data?
- This problem is called **model order selection**.



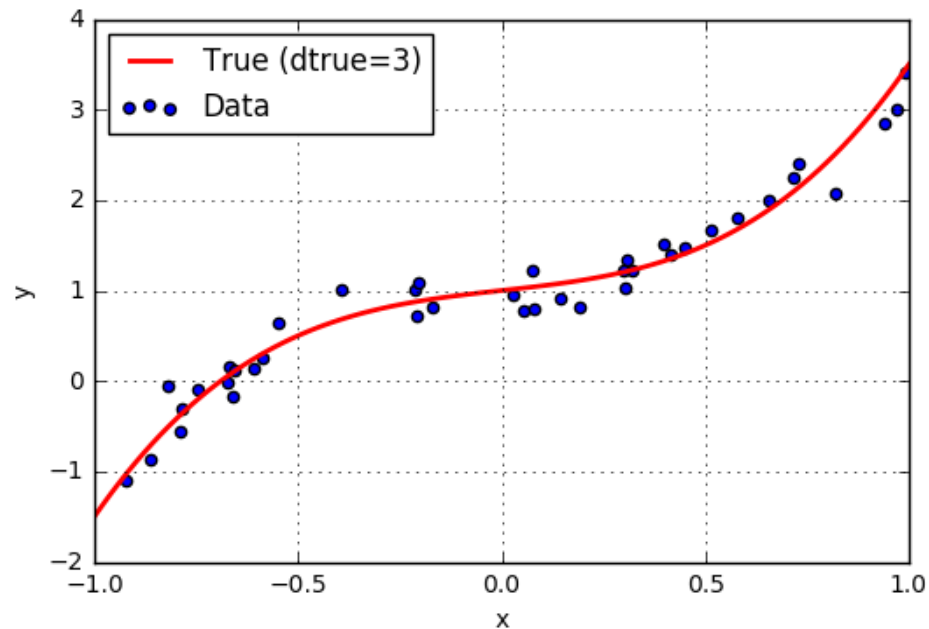
# Example Question

- You are given some data.
- Want to fit a model:  $y \approx f(x)$
- Decide to use a polynomial:  
 $f(x) = \beta_0 + \beta_1 x + \dots$
- What model order  $d$  should we use?
- Thoughts?



# Synthetic Data

- Previous example is synthetic data
  - $x_i$ : 40 samples uniform in  $[-1,1]$
  - $y = f(x) + \epsilon$ ,
    - $f(x) = \beta_0 + \beta_1 x + \dots + \beta_d x^d = \text{“true relation”}$
    - $d = 3, \epsilon \sim N(0, \sigma^2)$
- Synthetic data useful for analysis
  - Know “ground truth”
  - Can measure performance of various estimators



```
# Import useful polynomial library
import numpy.polynomial.polynomial as poly

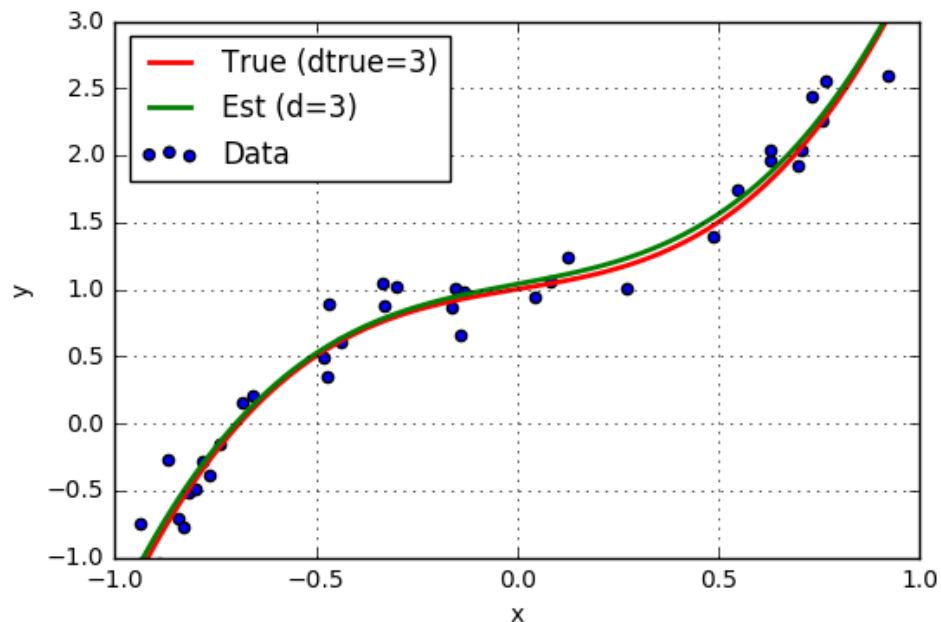
# True model parameters
beta = np.array([1,0.5,0,2]) # coefficients
wstd = 0.2 # noise
dtrue = len(beta)-1 # true poly degree

# Independent data
nsamp = 40
xdat = np.random.uniform(-1,1,nsamp)

# Polynomial
y0 = poly.polyval(xdat,beta)
ydat = y0 + np.random.normal(0,wstd,nsamp)
```

# Fitting with True Model Order

- Suppose true polynomial order,  $d=3$ , is known
- Use linear regression
  - `numpy.polynomial` package



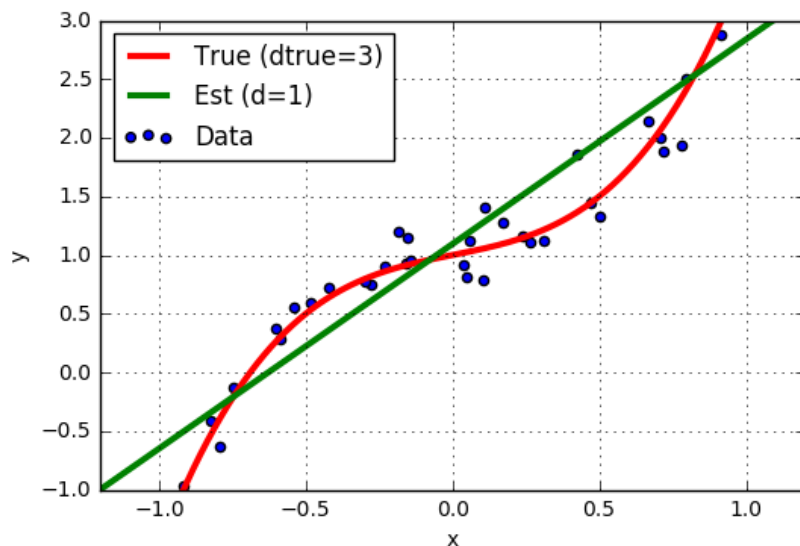
```
d = 3
beta_hat = poly.polyfit(xdat,ydat,d)

# Plot true and estimated function
xp = np.linspace(-1,1,100)
yp = poly.polyval(xp,beta)
yp_hat = poly.polyval(xp,beta_hat)
plt.xlim(-1,1)
plt.ylim(-1,3)
plt.plot(xp,yp,'r-',linewidth=2)
plt.plot(xp,yp_hat,'g-',linewidth=2)

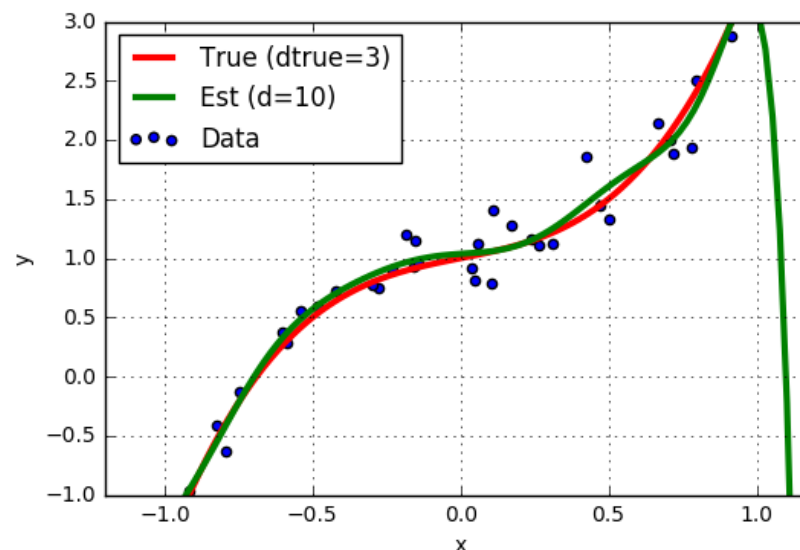
# Plot data
plt.scatter(xdat,ydat)
plt.legend(['True (dtrue=3)', 'Est (d=3)', 'Data'], loc='upper left')
plt.grid()
plt.xlabel('x')
plt.ylabel('y')
```

# But, True Model Order not Known

- Suppose we guess the wrong model order?

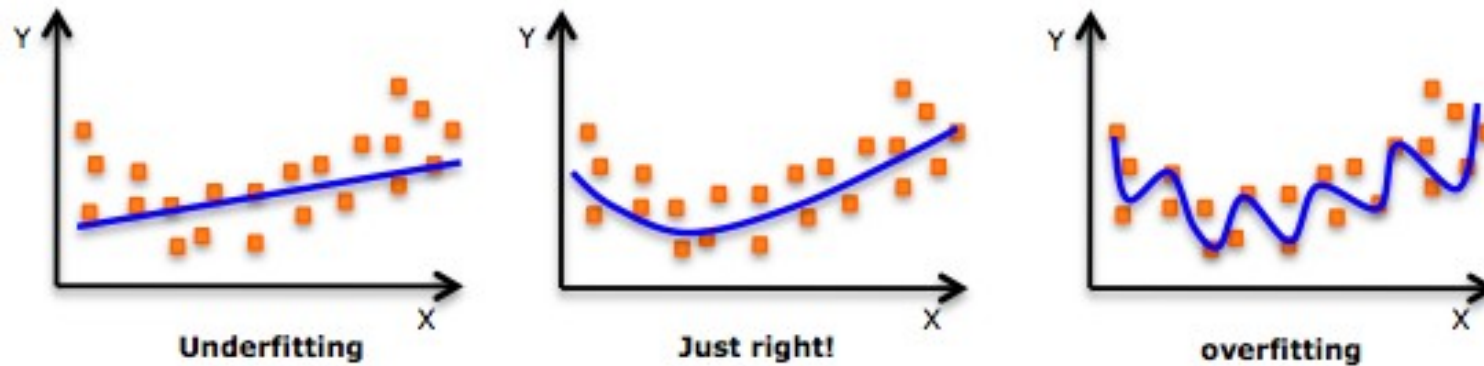


$d=1$  "Underfitting"



$d=10$  "Overfitting"

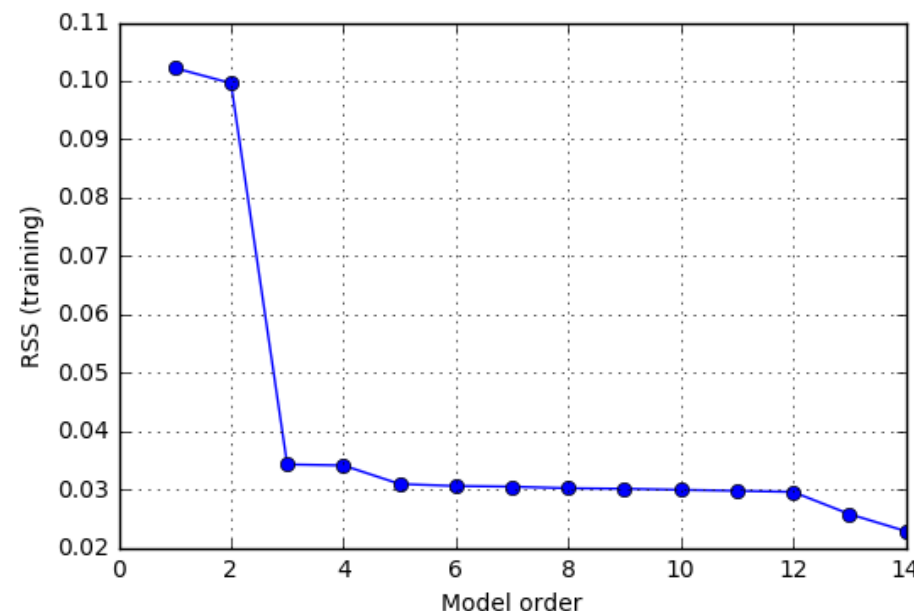
# How Can You Tell from Data?



- Is there a way to tell what is the correct model order to use?
- Must use the data. Do not have access to the true  $d$ ?
- What happens if we guess:
  - $d$  too big?
  - $d$  too small?

# Using RSS on Training Data?

- Simple (but bad) idea:
  - For each model order,  $d$ , find estimate  $\hat{\beta}$
  - Compute predicted values on training data
- $$\hat{y}_i = \hat{\beta}^T \mathbf{x}_i$$
- Compute RSS
$$RSS(d) = \sum_i (y_i - \hat{y}_i)^2$$
- Find  $d$  with lowest  $RSS$
- This doesn't work
  - $RSS(d)$  is always decreasing (Question: Why?)
  - Minimizing  $RSS(d)$  will pick  $d$  as large as possible
  - Leads to overfitting
- What went wrong?
- How do we do better?





# Model Class and True Function

- Analysis set-up:
  - Learning algorithm assumes a **model class**:  $\hat{y} = f(\mathbf{x}, \boldsymbol{\beta})$
  - But, data has **true** relation:  $y = f_0(x) + \epsilon, \epsilon \sim N(0, \sigma_\epsilon^2)$
- Will quantify three key effects:
  - Irreducible error
  - Under-modeling
  - Over-fitting

# Output Mean Squared Error

- To evaluate prediction error suppose we are given:
  - A parameter estimate  $\hat{\boldsymbol{\beta}}$  (computed from the learning algorithm)
  - A test point  $\mathbf{x}_{test}$
  - Test point is generally different from training samples.
- Predicted value:  $\hat{y} = f(\mathbf{x}_{test}, \hat{\boldsymbol{\beta}})$
- Actual value:  $y = f_0(\mathbf{x}_{test}) + \epsilon$
- Output mean squared error:
$$MSE_y(\mathbf{x}_{test}, \hat{\boldsymbol{\beta}}) := E[y - \hat{y}]^2$$
  - Expectation is over noise  $\epsilon$  on the test sample.

# Irreducible Error

- Rewrite output MSE:

$$MSE_y(\mathbf{x}_{test}, \hat{\boldsymbol{\beta}}) := E[y - \hat{y}]^2 = E[f_0(\mathbf{x}_{test}) + \epsilon - f(\mathbf{x}_{test}, \hat{\boldsymbol{\beta}})]^2$$

- Since noise on test sample is independent of  $\hat{\boldsymbol{\beta}}$  and  $\mathbf{x}_{test}$ :

$$MSE_y(\mathbf{x}_{test}, \hat{\boldsymbol{\beta}}) := [f_0(\mathbf{x}_{test}) - f(\mathbf{x}_{test}, \hat{\boldsymbol{\beta}})]^2 + E(\epsilon^2) = [f_0(\mathbf{x}_{test}) - f(\mathbf{x}_{test}, \hat{\boldsymbol{\beta}})]^2 + \sigma_\epsilon^2$$

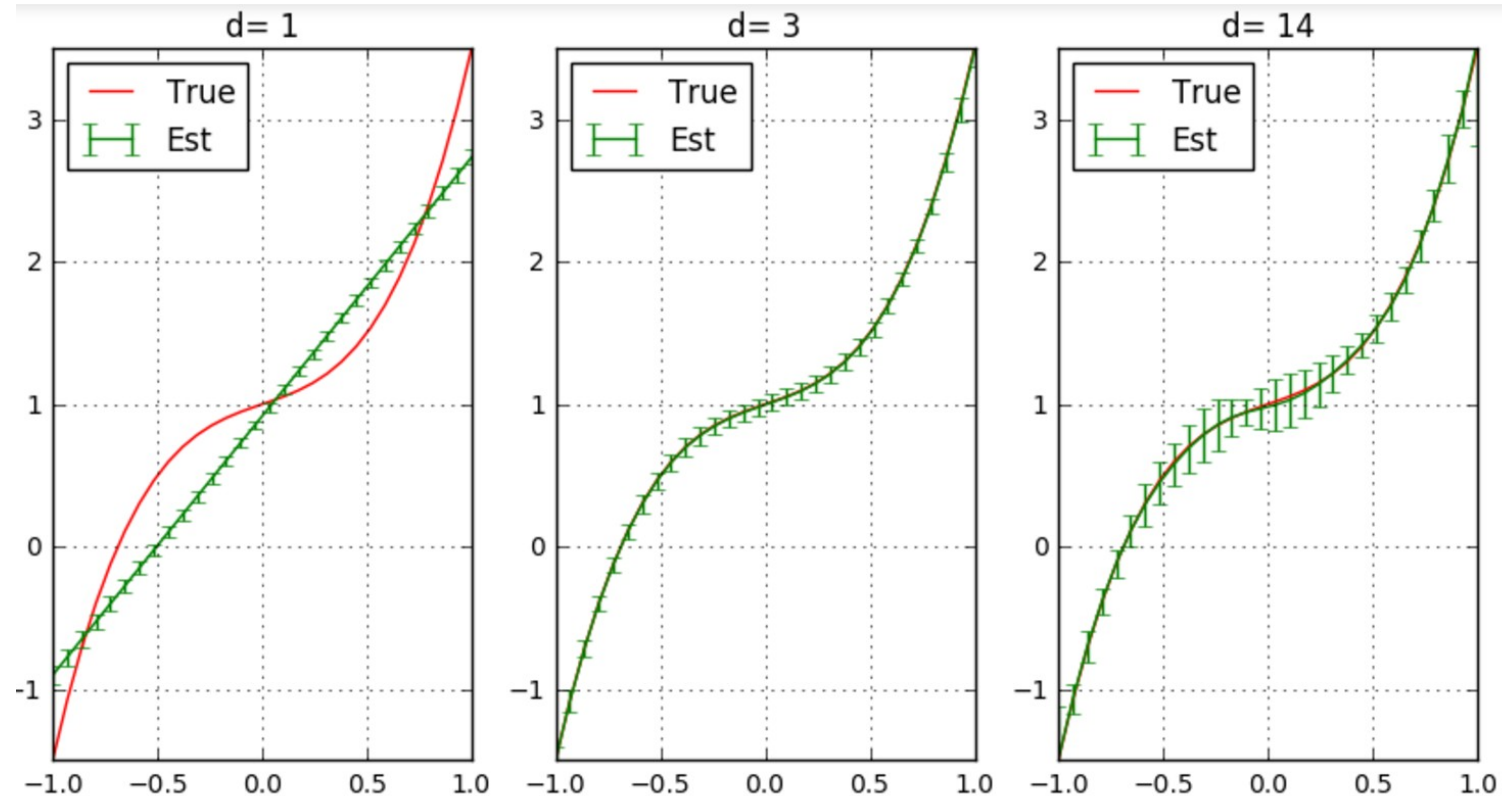
- Define **irreducible error**:  $\sigma_\epsilon^2$ 
  - Lower bound on  $MSE_y(\mathbf{x}_{test}, \hat{\boldsymbol{\beta}}) \geq \sigma_\epsilon^2$
  - Fundamental limit on ability to predict  $y$
  - Occurs since  $y$  is influenced by other factors than  $\mathbf{x}$

# Analysis with Noise (Advanced)

- - Now assume noise:  $y = f_0(\mathbf{x}) + \epsilon, \epsilon \sim N(0, \sigma_\epsilon^2)$
  - Get training data:  $(\mathbf{x}_i, y_i), i = 1, \dots, n$
  - Fit a parameter:
$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \sum_{i=1}^n (y_i - f(\mathbf{x}_i, \boldsymbol{\beta}))^2$$
    - $\hat{\boldsymbol{\beta}}$  will be random.
    - Depends on particular noise realization.
  - Take a new test point  $\mathbf{x}_{test}$  (not random)
  - Compute mean and variance of estimated function  $f(\mathbf{x}_{test}, \hat{\boldsymbol{\beta}})$
  - Define:
    - **Bias**: Difference of true function from mean estimate
    - **Variance**: Variance of estimate around its mean

# Bias and Variance Illustrated

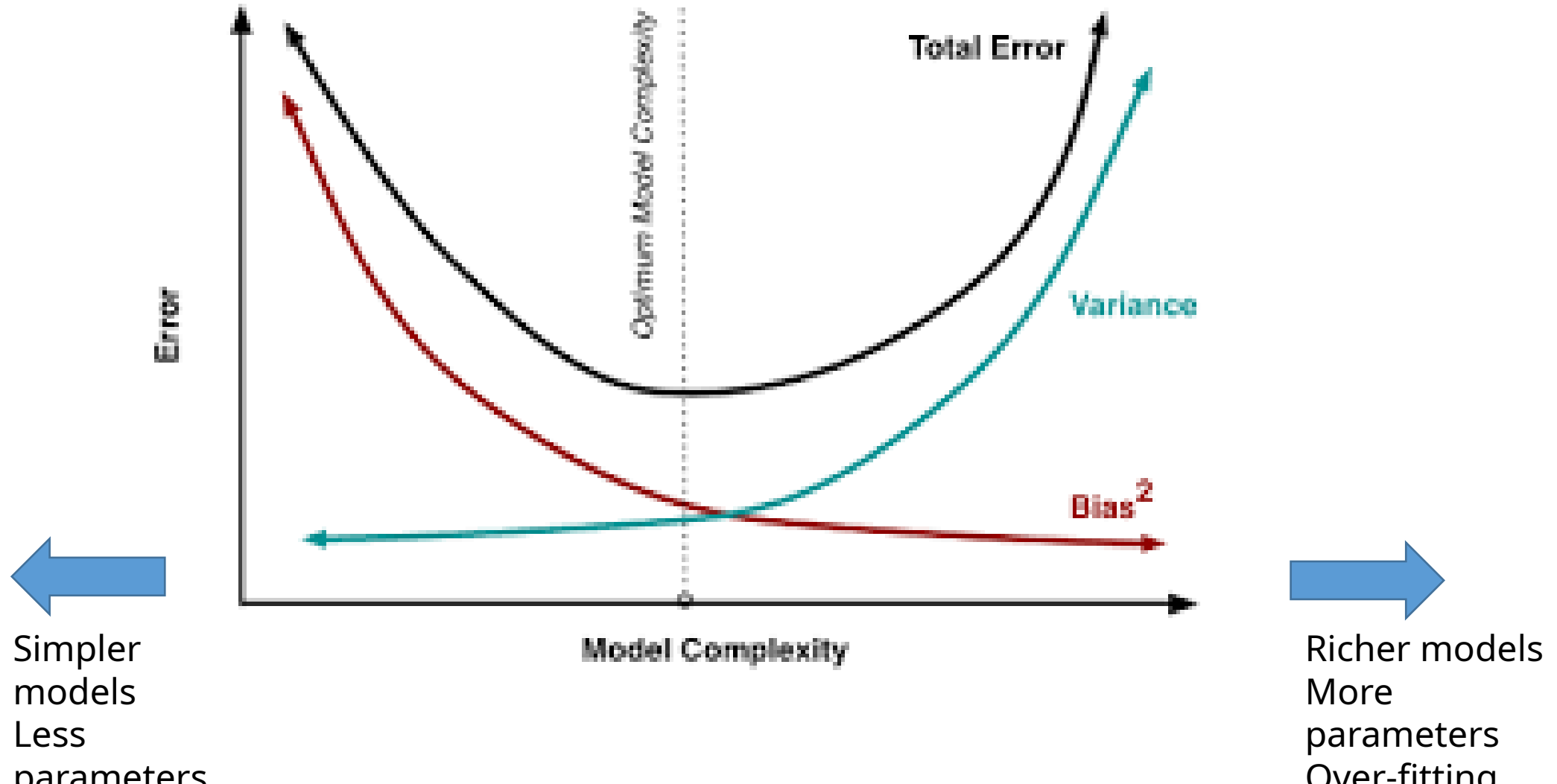
- Polynomial ex
- Mean and std dev of estimated functions
- 100 trials



Low  
variance,  
High bias

High  
variance,  
Zero bias

# Bias-Variance Tradeoff



# Cross Validation

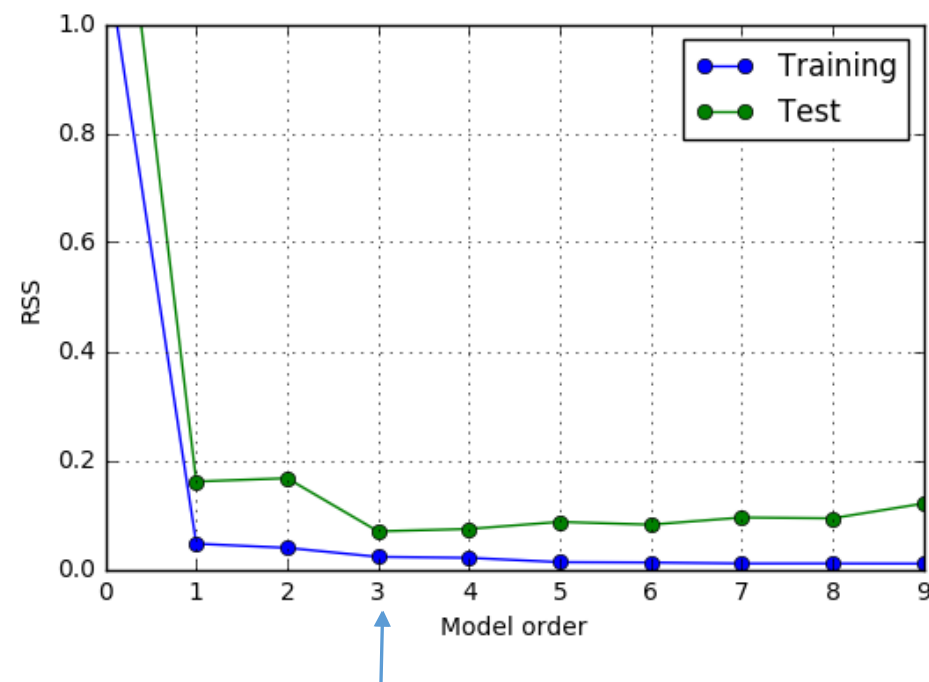
- Concept: Need to test fit on data independent of training data
- Divide data into two sets:
  - $N_{train}$  training samples,  $N_{valid}$  validation samples
- For each model order,  $p$ , learn parameters  $\hat{\beta}$  from training samples
- Measure RSS on validation samples.

$$RSS_{test}(p) = \sum_{i \in \text{valid}} (\hat{y}_i - y_i)^2$$

- Select model order  $p$  that minimizes  $RSS_{valid}(p)$

# Finding the Model Order

- Estimated optimal model order = 3



RSS test minimized at  $d = 3$   
RSS training always decreases

```
dtest = np.array(range(0,10))
RSStest = []
RSStr = []
for d in dtest:

    # Fit data
    beta_hat = poly.polyfit(xtr,ytr,d)

    # Measure RSS on training data
    # This is not necessary, but we do it just to show the training error
    yhat = poly.polyval(xtr,beta_hat)
    RSSd = np.mean((yhat-ytr)**2)
    RSStr.append(RSSd)

    # Measure RSS on test data
    yhat = poly.polyval(xts,beta_hat)
    RSSd = np.mean((yhat-yts)**2)
    RSStest.append(RSSd)

plt.plot(dtest,RSStr,'bo-')
plt.plot(dtest,RSStest,'go-')
plt.xlabel('Model order')
plt.ylabel('RSS')
plt.grid()
plt.ylim(0,1)
plt.legend(['Training','Test'],loc='upper right')
```

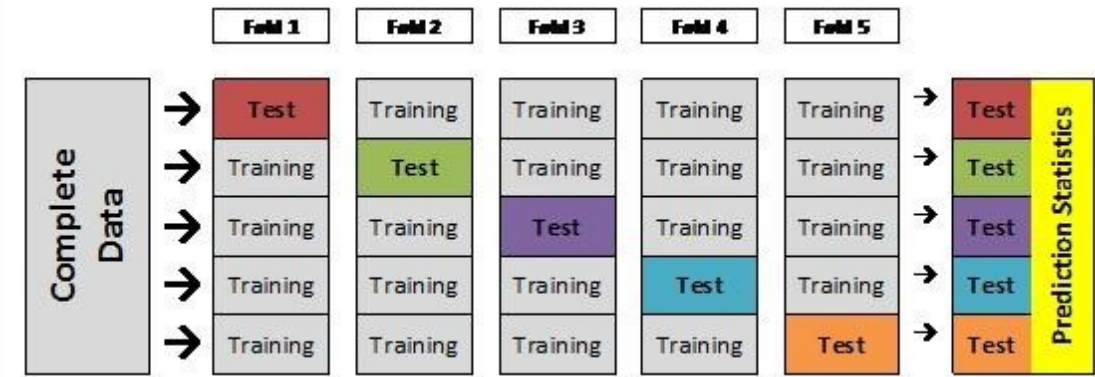


# Problems with Simple Train/Test Split

- Test error could vary significantly depending on samples selected
- Only use limited number of samples for training
- Problems particularly bad for data with limited number of samples

# K-Fold Cross Validation

- $K$ -fold cross validation
  - Divide data into  $K$  parts
  - Use  $K - 1$  parts for training. Use remaining for test.
  - Average over the  $K$  test choices
  - More accurate, but requires  $K$  fits of parameters
- Leave one out cross validation (LOOCV)
  - Take  $K = N$  so one sample is left out.
  - Most accurate, but requires  $N$  model fittings



# Polynomial Example

- Use sklearn Kfold object
- Loop
  - Outer loop: Over K folds
  - Inner loop: Over model ord
  - Measure test error in each f
  - Can be time-consuming

```
# Create a k-fold object
nfold = 20
kf = sklearn.model_selection.KFold(n_splits=nfold,shuffle=True)

# Model orders to be tested
dtest = np.arange(0,10)
nd = len(dtest)

# Loop over the folds
RSSs = np.zeros((nd,nfold))
for isplit, Ind in enumerate(kf.split(xdat)):

    # Get the training data in the split
    Itr, Its = Ind
    xtr = xdat[Itr]
    ytr = ydat[Itr]
    xts = xdat[Its]
    yts = ydat[Its]

    for it, d in enumerate(dtest):

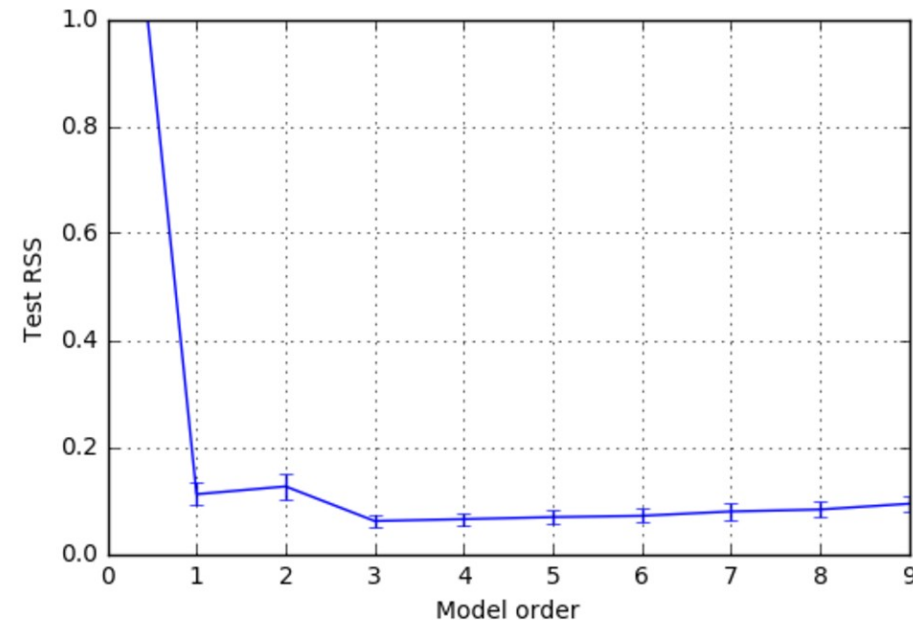
        # Fit data on training data
        beta_hat = poly.polyfit(xtr,ytr,d)

        # Measure RSS on test data
        yhat = poly.polyval(xts,beta_hat)
        RSSs[it,isplit] = np.mean((yhat-yts)**2)
```

# Polynomial Example CV Results

- For each model order  $d$ 
  - Compute mean test RSS
  - Compute std error (SE) of test RSS
  - $SE = \text{std dev} / \sqrt{K - 1}$
  - Mean and SE computed over the  $K$  folds
- Simple model selection
  - Select  $d$  with lowest mean test RSS
- For this example
  - Estimate model order = 3

```
RSS_mean = np.mean(RSSsts,axis=1)
RSS_std = np.std(RSSsts,axis=1) / np.sqrt(nfold-1)
plt.errorbar(dtest, RSS_mean, yerr=RSS_std, fmt='-')
plt.ylim(0,1)
plt.xlabel('Model order')
plt.ylabel('Test RSS')
plt.grid()
```

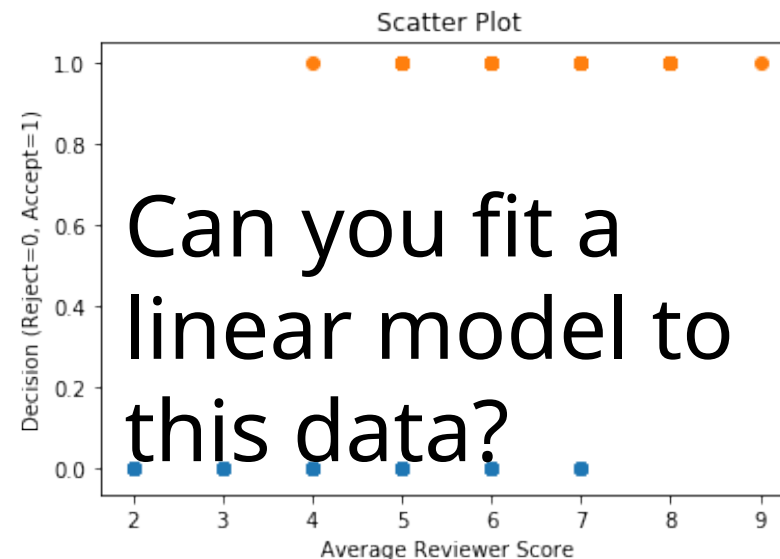


# Binary Classification

## Binary Classification

### Task (T):

- Simplest example where  $x \in \mathbb{R}$  and  $y \in \{0, 1\}$



- Dataset of ICLR'18 review scores vs. accept/reject decisions

TL;DR	_bibtex	abstract	authorids	authors	conf_1	conf_2	conf_3	decision	review	review_1	review_2	review_3	title
None	@article{\nsharma2018hyperedge2vec,\ntitle={H...	Data structured in form of overlapping or non-...	[sharm170@umn.edu, srjoty@ntu.edu.sg, himanshu...	[Ankit Sharma, Shafiq Joty, Himanshu Kharkwal,...	3.0	3.0	4.0	Reject	5.000000	5.0	5.0	5.0	Hyperedge2vec: Distributed Representations for...
Query-based black-box attacks on deep neural n...	@article{\nnitin2018exploring,\ntitle={Explori...	Existing black-box attacks on deep neural netw...	[abhagoji@princeton.edu, _w@eecs.berkeley.edu,...	[Arjun Nitin Bhagoji, Warren He, Bo Li, Dawn S...	4.0	3.0	4.0	Reject	6.000000	6.0	6.0	7.0	Exploring the Space of Black-box Attacks on De...
A theory and algorithmic framework for predict...	@article{\nd.2018learning,\ntitle={Learning We...	Predictive models that generalize well under d...	[fredrikj@mit.edu, kallus@cornell.edu, urish22...	[Fredrik D. Johansson, Nathan Kallus, Uri Shal...	3.0	3.0	4.0	Reject	6.666667	5.0	8.0	7.0	Learning Weighted Representations for Generali...
We prove that DNN is a recursively approximate...	@article{\nzheng2018understanding,\ntitle={Und...	Deep learning achieves remarkable generalizati...	[zhenggh@mail.ustc.edu.cn, jtsang@bjtu.edu.cn,...	[Guanhua Zheng, Jitao Sang, Changsheng Xu]	3.0	3.0	2.0	Reject	3.666667	2.0	3.0	6.0	Understanding Deep Learning Generalization by

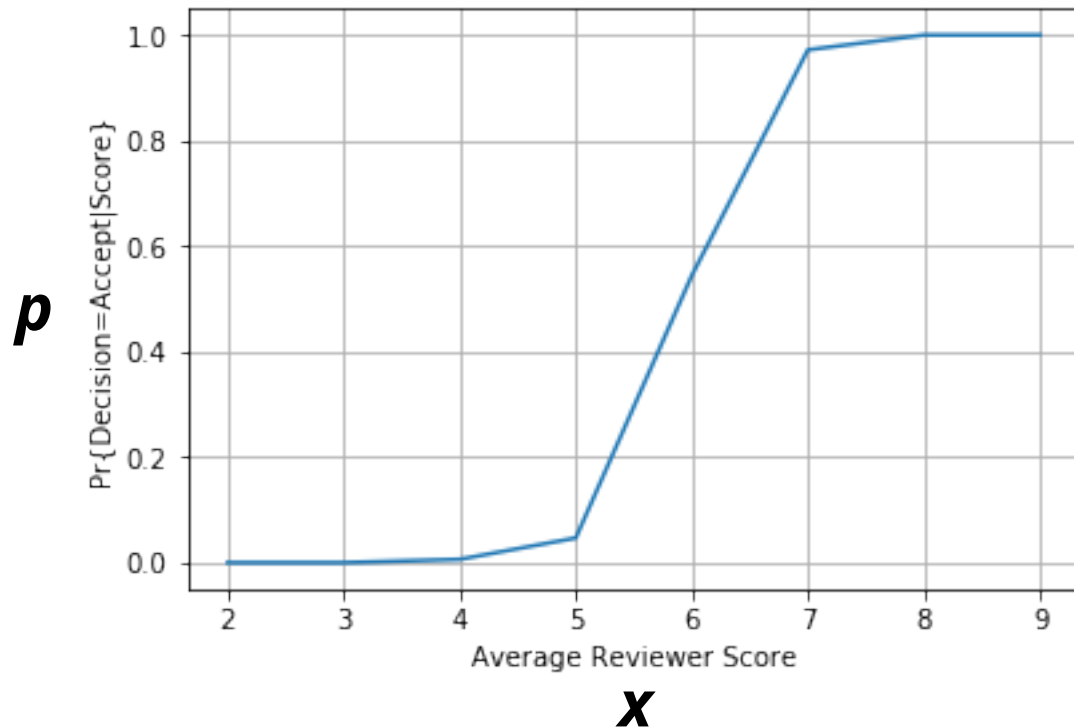
# Logistic Regression

## Binary Classification

### Task (T):

Instead, let's compute and plot  $\Pr\{y = 1 \mid x\}$

$\Pr\{\text{Decision}=\text{Accept} \mid \text{Score}\}$   
↓  
 $\Pr\{y = 1 \mid x\}$



- Idea: Linear regression to fit  $p$  as a function of  $x$

$$p = \beta_1 x + \beta_0$$

- Is this a good idea?
  - Probability  $p$  is always bounded between  $[0,1]$

# Logistic Regression

## Binary Classification

### Task (T):

Consider the following function:

"Logits"  
Function

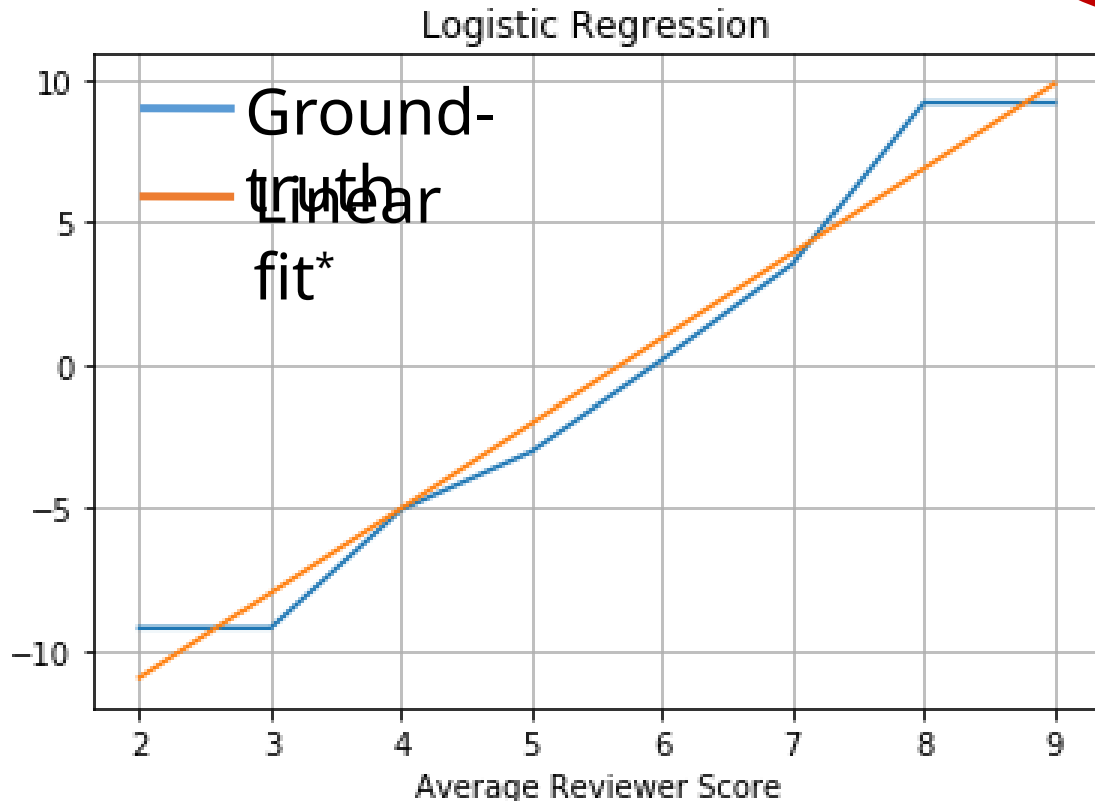
$$g = \log\left(\frac{p}{1-p}\right)$$

- What is the range of  $g$ ?

$$g \in [-\infty, \infty]$$

- Logistic Regression: **fit logits function using a linear model!**

$$g = \log\left(\frac{p}{1-p}\right) = \beta_1 x + \beta_0$$



**x**

Note: the linear fit is illustrative only. How to determine the best linear fit will be

# Logistic Regression

$$g = \log\left(\frac{p}{1-p}\right) = \beta_1 x + \beta_0 \quad \Rightarrow \quad \begin{array}{c} \text{Pr}\{\text{Decision}=\text{Accept} | \\ \text{Score}\} \\ p = \frac{1}{1 + e^{-(\beta_1 x + \beta_0)}} \end{array}$$

- What is  $\text{Pr}\{\text{Decision}=\text{Reject} | \text{Score}\}$

$$1 - p = \frac{e^{-(\beta_1 x + \beta_0)}}{1 + e^{-(\beta_1 x + \beta_0)}}$$

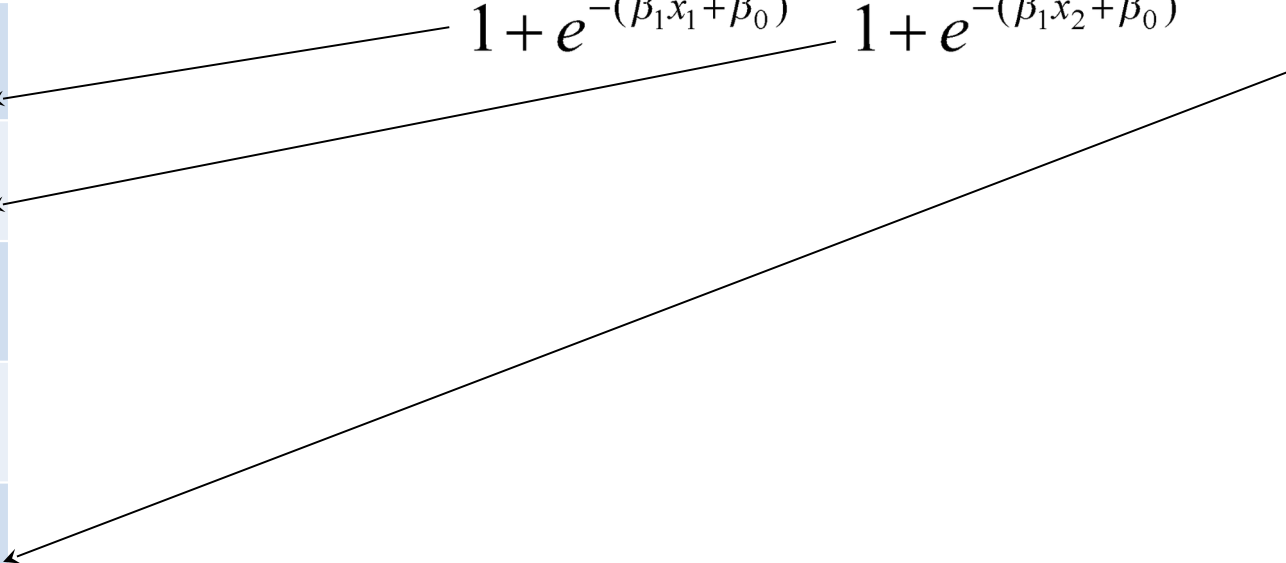
How do we find the model parameters  $\beta_1$  and  $\beta_0$ ?



# Model Estimation

- We will use an approach referred to as **Maximum Likelihood Estimation** (MLE)
  - Let's assume that the model(i.e.,  $\beta_1$  and  $\beta_0$ ) is magically known. Consider the training dataset below. What is the likelihood that the dataset came from our model?

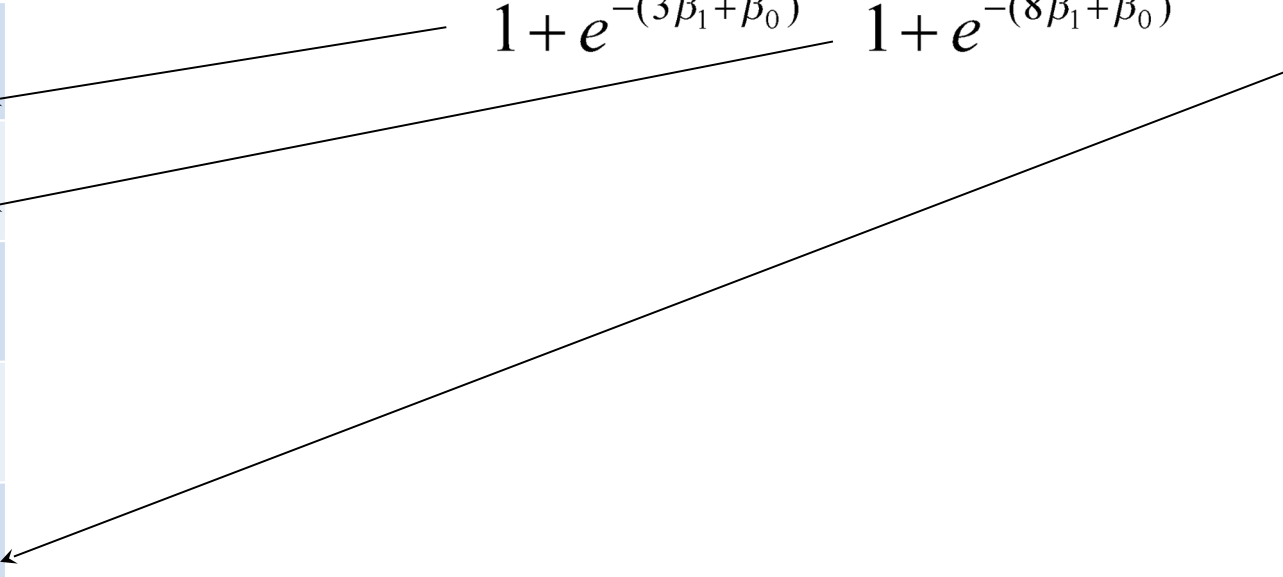
#	X	Y
1	$x_1 = 3$	$y_1 = 0$
2	$x_2 = 8$	$y_2 = 1$
..		
..		
N	$x_N = 6$	$y_N = 1$

$$\text{Likelihood} = \frac{e^{-(\beta_1 x_1 + \beta_0)}}{1 + e^{-(\beta_1 x_1 + \beta_0)}} * \frac{1}{1 + e^{-(\beta_1 x_2 + \beta_0)}} * \dots * \frac{1}{1 + e^{-(\beta_1 x_N + \beta_0)}}$$


# Model Estimation

- We will use an approach referred to as **Maximum Likelihood Estimation** (MLE)
  - Let's assume that the model(i.e.,  $\beta_1$  and  $\beta_0$ ) is magically known. Consider the training dataset below. What is the likelihood that the dataset came from our model?

#	X	Y
1	$x_1 = 3$	$y_1 = 0$
2	$x_2 = 8$	$y_2 = 1$
..		
..		
N	$x_N = 6$	$y_N = 1$

$$\text{Likelihood} = \frac{e^{-(3\beta_1 + \beta_0)}}{1 + e^{-(3\beta_1 + \beta_0)}} * \frac{1}{1 + e^{-(8\beta_1 + \beta_0)}} * \dots * \frac{1}{1 + e^{-(6\beta_1 + \beta_0)}}$$


# Model Estimation

- We will use an approach referred to as **Maximum Likelihood Estimation** (MLE)
  - Let's assume that the model(i.e.,  $\beta_1$  and  $\beta_0$ ) is magically known. Consider the training dataset below. What is the likelihood that the dataset came from our model?

#	X	Y
1	$x_1 = 3$	$y_1 = 0$
2	$x_2 = 8$	$y_2 = 1$
..		
..		
N	$x_N = 6$	$y_N = 1$

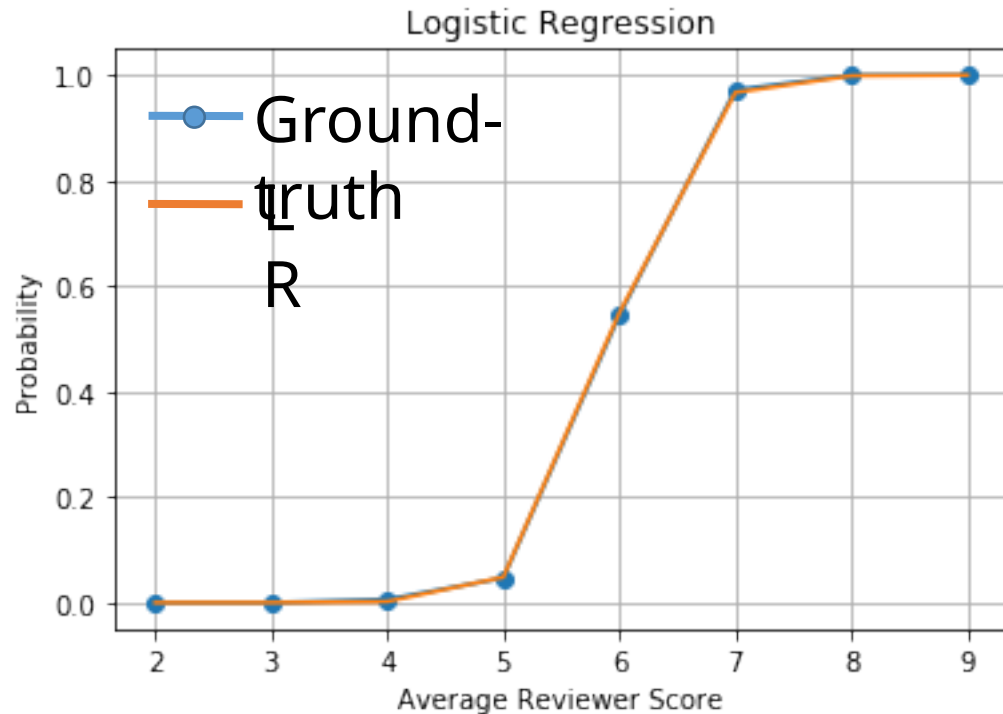
Log-Likelihood =  $\log\left(\frac{e^{-(3\beta_1+\beta_0)}}{1+e^{-(3\beta_1+\beta_0)}}\right) + \log\left(\frac{1}{1+e^{-(8\beta_1+\beta_0)}}\right) + \dots \log\left(\frac{1}{1+e^{-(6\beta_1+\beta_0)}}\right)$

$g(\beta_1, \beta_0)$  Function of model parameters

only  
Find  $\beta_1$  and  $\beta_0$  that  
maximize g

(or minimize the "loss" -g)  
 $Loss(\beta_1, \beta_0) = -g(\beta_1, \beta_0)$

# We Won't Worry About How (Phew!)



```
from sklearn import linear_model

#Instantiate an LR object
logreg = sklearn.linear_model.LogisticRegression(C=1e5);

#Recall: your training data must have a column of ones for the constant term
xd = np.ones((numPapers,2));
xd[:,0] = np.append(rscores,ascores)

yd = np.append(rlabels,alabels);

logreg.fit(xd,yd);

#Plot Pr{Accept|Score}
rv = np.ones((len(revRange),2));
rv[:,0] = revRange;
prpredict=logreg.predict_proba(rv)
```

**From regression to classification: if probability of Accept > 0.5, then output Accept.**

# Logistic Regression: Multi-Variate Case

## UCI Spam

## Dataset:

<https://archive.ics.uci.edu/ml/datasets/>

### Spambase

#### Attribute Information

The last column of 'spambase.data' denotes whether the e-mail was considered spam (1) or not (0), i.e. unsolicited commercial e-mail. Most of the attributes indicate whether a particular word or character was frequently occurring in the e-mail. The run-length attributes (55-57) measure the length of sequences of consecutive capital letters. For the statistical measures of each attribute, see the end of this file. Here are the definitions of the attributes:

48 continuous real [0,100] attributes of type word\_freq\_WORD

= percentage of words in the e-mail that match WORD, i.e.  $100 * (\text{number of times the WORD appears in the e-mail}) / \text{total number of words in e-mail}$ . A "word" in this case is any string of alphanumeric characters bounded by non-alphanumeric characters or end-of-string.

6 continuous real [0,100] attributes of type char\_freq\_CHAR

= percentage of characters in the e-mail that match CHAR, i.e.  $100 * (\text{number of CHAR occurrences}) / \text{total characters in e-mail}$

1 continuous real [1,...] attribute of type capital\_run\_length\_average

= average length of uninterrupted sequences of capital letters

1 continuous integer [1,...] attribute of type capital\_run\_length\_longest

= length of longest uninterrupted sequence of capital letters

1 continuous integer [1,...] attribute of type capital\_run\_length\_total

= sum of length of uninterrupted sequences of capital letters

= total number of capital letters in the e-mail

1 nominal {0,1} class attribute of type spam

= denotes whether the e-mail was considered spam (1) or not (0), i.e. unsolicited commercial e-mail.

- **57 Real or integer valued features**

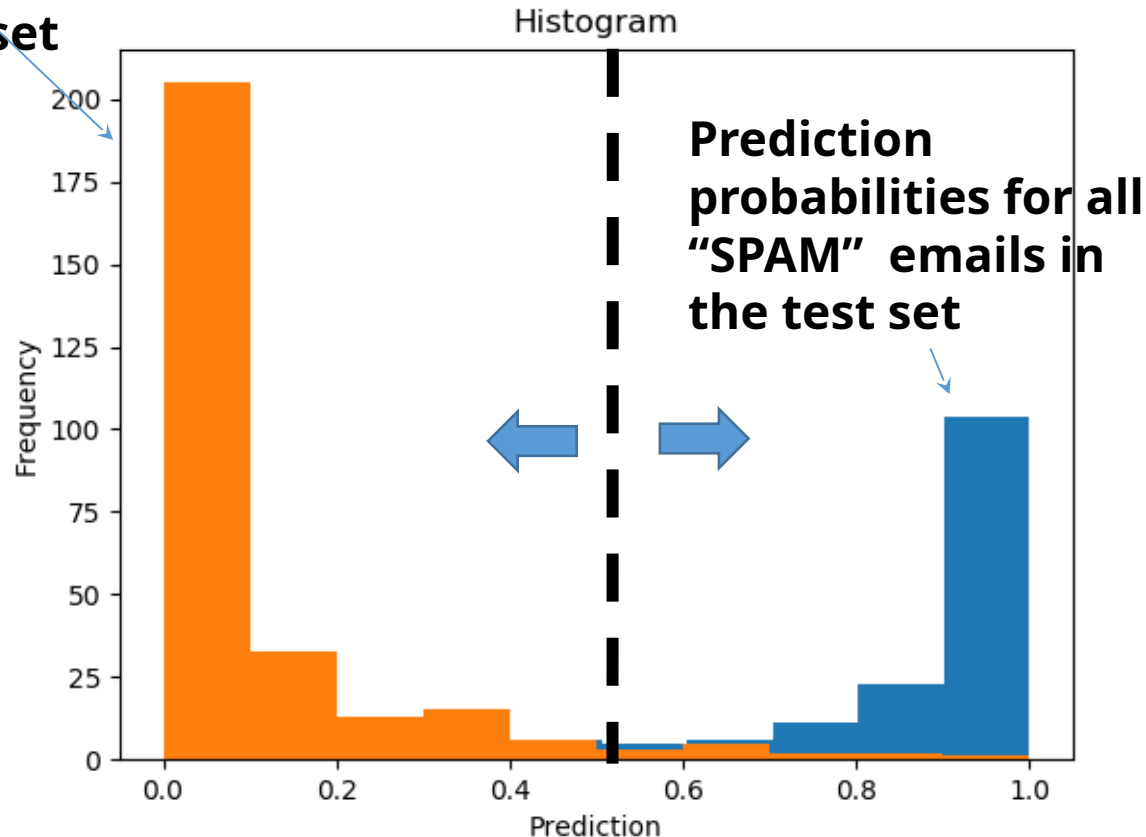
- **Binary output class**

$$p_{spam} = \frac{1}{1 + e^{-(\sum_{i=1}^M \beta_i x_i + \beta_0)}}$$

# LR on Spam Database: Results

90% of samples used for training, remaining 10%  
used for test

**Prediction probabilities  
for all "SPAM" emails in  
the test set**



```
#Instantiate an LR object
logreg = sklearn.linear_model.LogisticRegression(C=1e5);

#Recall: your training data must have a column of ones for the constant term
xd = np.ones((numPapers,2));
xd[:,0] = np.append(rscores,ascores)

yd = np.append(rlabels,alabels);

logreg.fit(xd,yd);

#Plot Pr{Accept|Score}
rv = np.ones((len(revRange),2));
rv[:,0] = revRange;
prpredict=logreg.predict_proba(rv)
```

Which emails are mis-  
predicted?  
Accuracy on test set:  
~92%

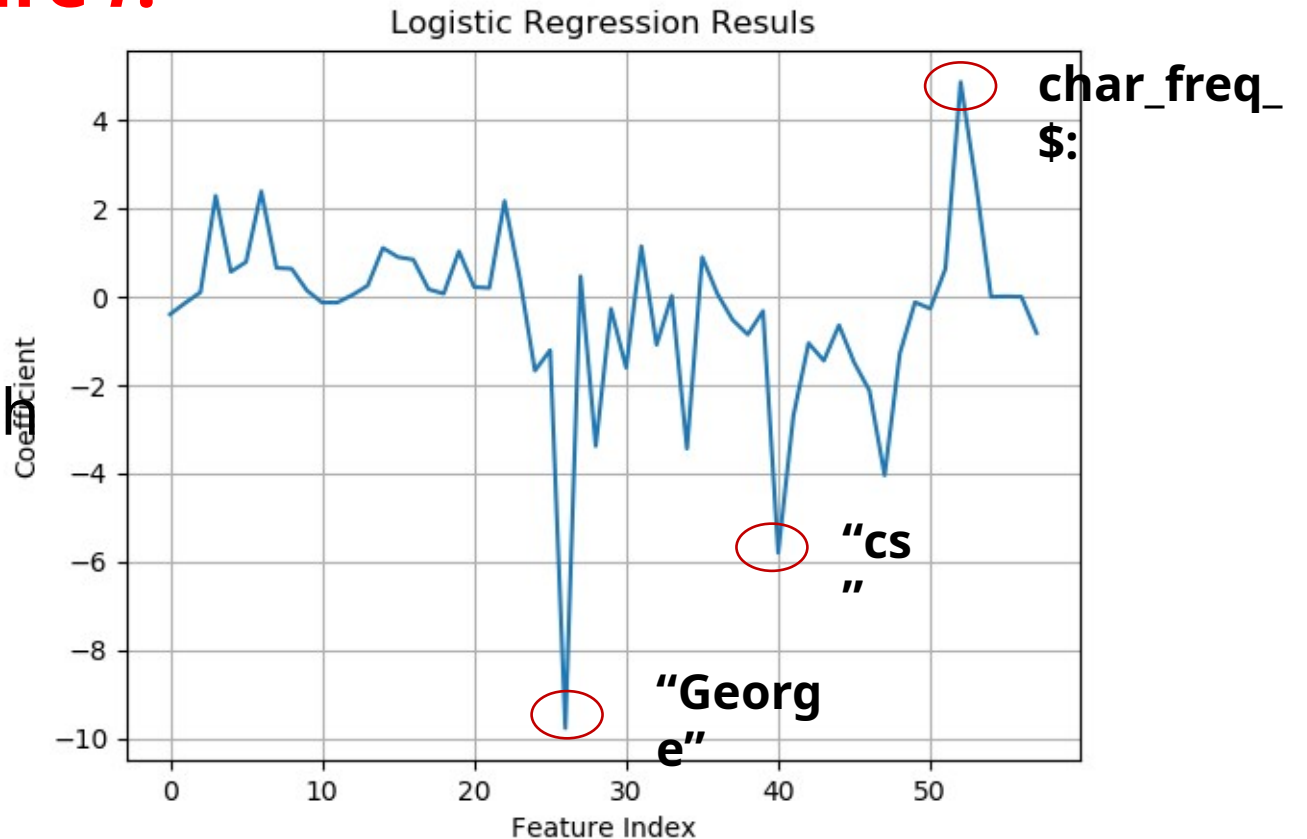
# Which Features Matter?

Our  
Model:

$$p_{spam} = \frac{1}{1 + e^{-(\sum_{i=1}^M \beta_i x_i + \beta_0)}}$$

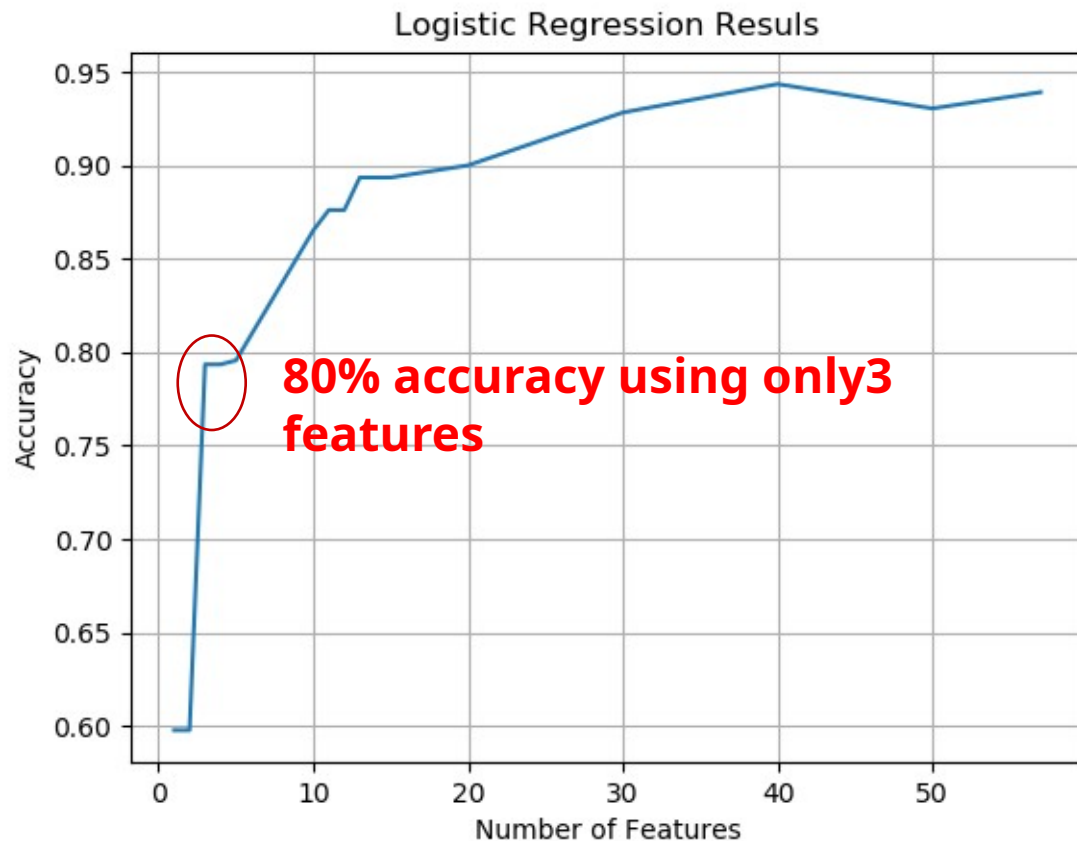
Reasonable hypothesis: features with larger absolute values of  $\beta$  matter more.

What does  $\beta_i=0$  imply about feature  $i$ ?



# Feature Selection

Retrain and predict using only the top-k features



Can we explicitly train the parameters so as to prioritize a “sparser” model?

**Why?** Low model complexity prevents overfitting

Recall that during training we were seeking to minimize:

$$\hat{\beta} = \min_{\beta} Loss(\beta)$$

**How should this objective function change?**



# Regularization

**$L_p$  Norm of a vector  $x$**

$$\|x\|_p = (\sum |x_i|^p)^{1/p}$$

$p$	$L_p$ Norm	Interpretation
0	$\ x\ _0 = (\sum  x_i ^0)^{1/0}$	Number of Non-zero Entries
1	$\ x\ _1 = (\sum  x_i )$	Sum of absolute values
2	$\ x\ _2 = (\sum  x_i ^2)^{0.5}$	Root mean square
$\infty$	$\ x\ _\infty = (\sum  x_i ^\infty)^0$	Max. value

**“Regularized”  
loss**

$$\hat{\beta} = \min_{\beta} \{Loss(\beta) + c\|\beta\|_0\}$$

$c$  controls the relative importance of the regularization penalty

# Regularization In Practice

**L0  
Regularization**

$$\hat{\beta} = \min_{\beta} \{Loss(\beta) + c\|\beta\|_0\}$$

**Hard  
“combinatorial”  
optimization  
problem!**

**Instead, the following regularization functions are commonly used:**

**L1  
Regularization**

$$\hat{\beta} = \min_{\beta} \{Loss(\beta) + c\|\beta\|_1\}$$

**(LASSO)**

**L2  
Regularization**

$$\hat{\beta} = \min_{\beta} \{Loss(\beta) + c\|\beta\|_2\}$$

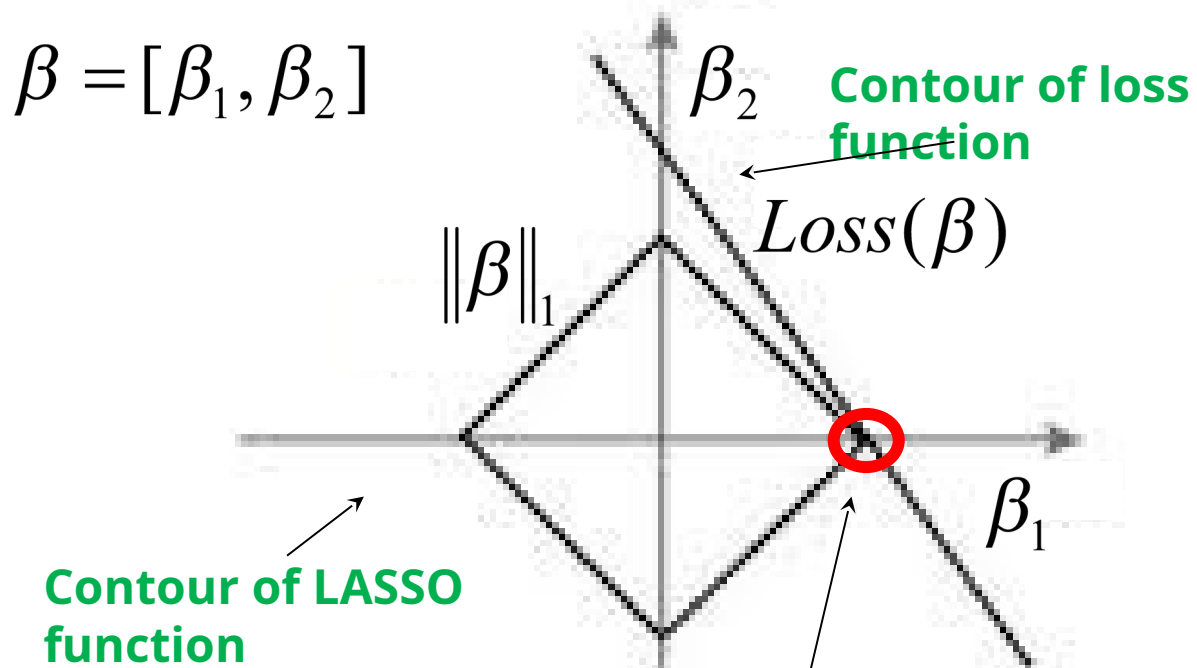
**(Ridge)**

We are penalizing “large” coefficients.

But why?

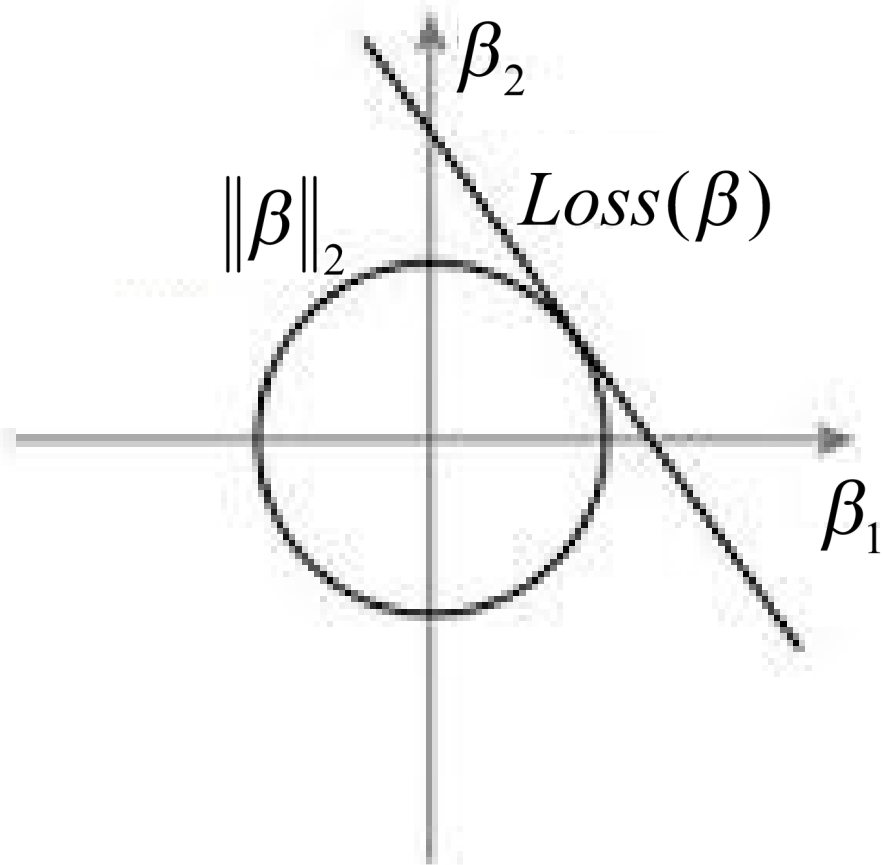
# LASSO and Ridge Regularization

**A** L1 regularization



**LASSO prefers  
sparse  
solutions!**

**B** L2 regularization



# Regularization for Spam Classification

Logistic Regression (aka logit, MaxEnt) classifier.

In the multiclass case, the training algorithm uses the one-vs-rest (OvR) scheme if the 'multi\_class' option is set to 'ovr', and uses the cross-entropy loss if the 'multi\_class' option is set to 'multinomial'. (Currently the 'multinomial' option is supported only by the 'lbfgs', 'sag' and 'newton-cg' solvers.)

This class implements regularized logistic regression using the 'liblinear' library, 'newton-cg', 'sag' and 'lbfgs' solvers. It can handle both dense and sparse input. Use C-ordered arrays or CSR matrices containing 64-bit floats for optimal performance; any other input format will be converted (and copied).

The 'newton-cg', 'sag', and 'lbfgs' solvers support only L2 regularization with primal formulation. The 'liblinear' solver supports both L1 and L2 regularization, with a dual formulation only for the L2 penalty.

[Read more in the User Guide.](#)

**Parameters:** `penalty` : str, 'l1' or 'l2', default: 'l2'

Used to specify the norm used in the penalization. The 'newton-cg', 'sag' and 'lbfgs' solvers support only l2 penalties.

*New in version 0.19:* l1 penalty with SAGA solver (allowing 'multinomial' + L1)

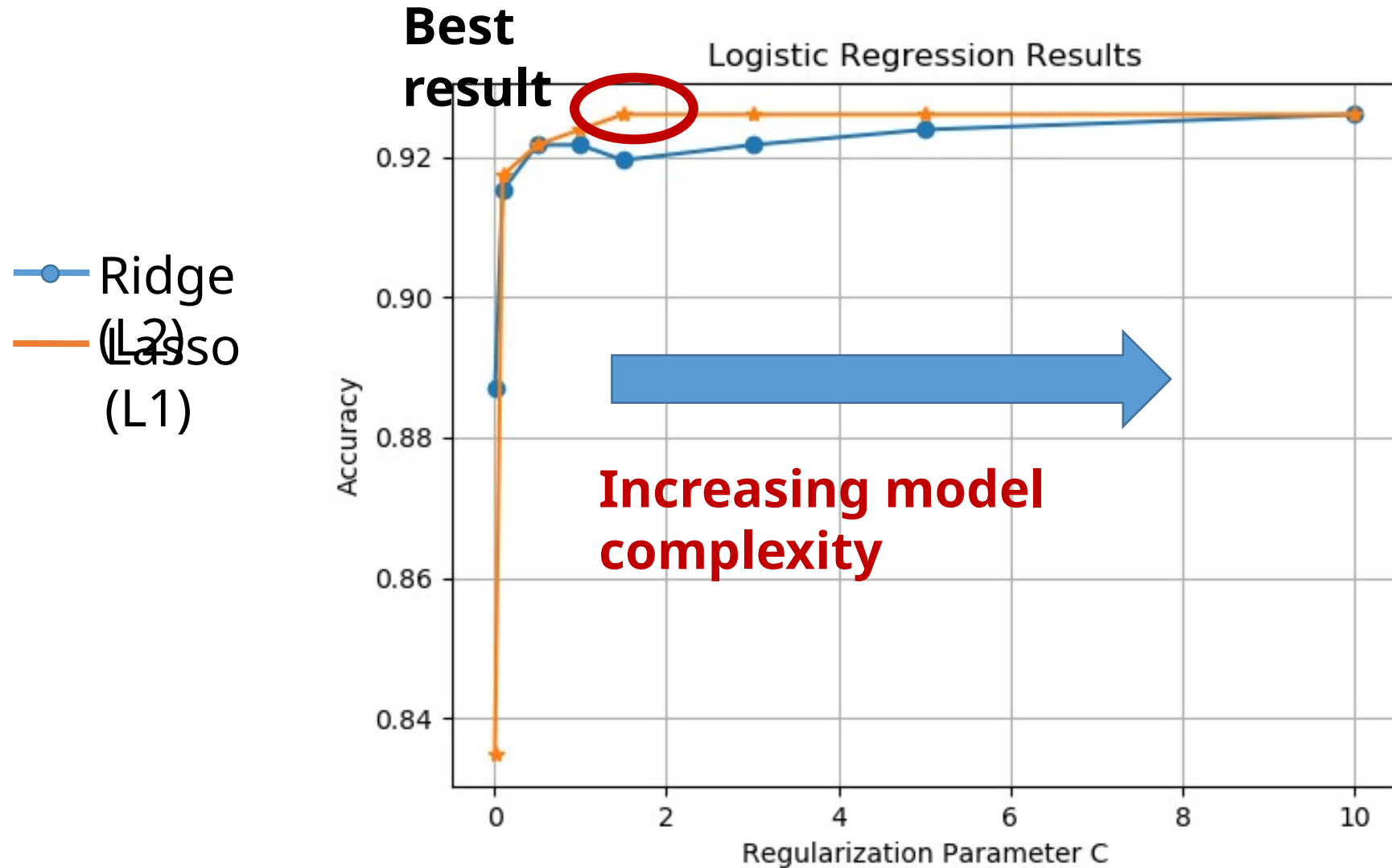
`C` : float, default: 1.0

Inverse of regularization strength; must be a positive float. Like in support vector machines, smaller values specify stronger regularization.

Which regularization function to use?

How should we select  $c$ ?

# Impact of $C$



# Errors in Binary Classification

- Two types of errors:
  - Type I error (False positive / false alarm): Decide  $\hat{y} = 1$  when  $y = 0$
  - Type II error (False negative / missed detection): Decide  $\hat{y} = 0$  when  $y = 1$
- Implication of these errors may be different
  - Think of breast cancer diagnosis
- Accuracy of classifier can be measured by:
  - $TPR = P(\hat{y} = 1|y = 1)$
  - $FPR = P(\hat{y} = 1|y = 0)$

predicted→ real↓	Class_pos	Class_neg
Class_pos	TP	FN
Class_neg	FP	TN

$$TPR \text{ (sensitivity)} = \frac{TP}{TP + FN}$$

$$FPR \text{ (1-specificity)} = \frac{FP}{TN + FP}$$

# ROC Curve

```
from sklearn import metrics
yprob = logreg.predict_log_proba(Xtr)
fpr, tpr, thresholds = metrics.roc_curve(ytr,yprob[:,1])

plt.loglog(fpr,1-tpr)
plt.grid()
plt.xlabel('FPR')
plt.ylabel('TPR')
```

- Varying threshold obtains a set of classifier
- Trades off FPR and TPR
- Can visualize with ROC curve
  - Receiver operating curve
  - Term from digital communications

