

Lecture 2: ML Basics + Spam Filtering

Siddharth Garg
sg175@nyu.edu

Classification

Task (T):

- Emails $x \in$ all possible emails and $y \in \{spam, non_spam\}$ find

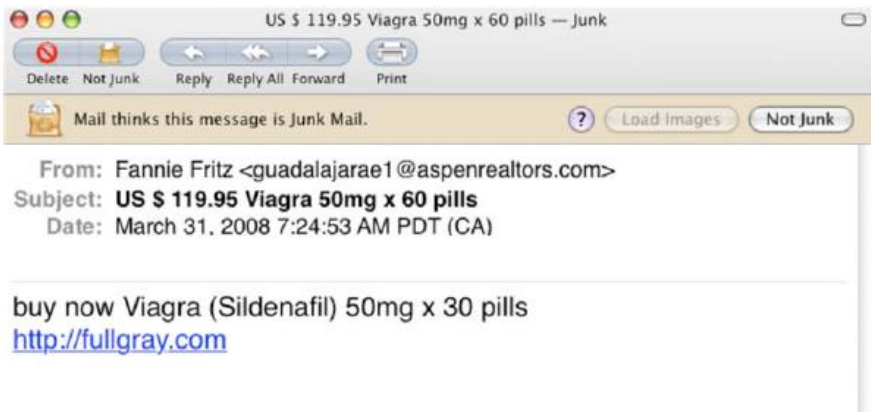
$$f: x \rightarrow y$$

Experience (E):

- A "training dataset" a emails marked as "spam" or "non_spam"

Performance (P):

- Spam detection accuracy



SPAM

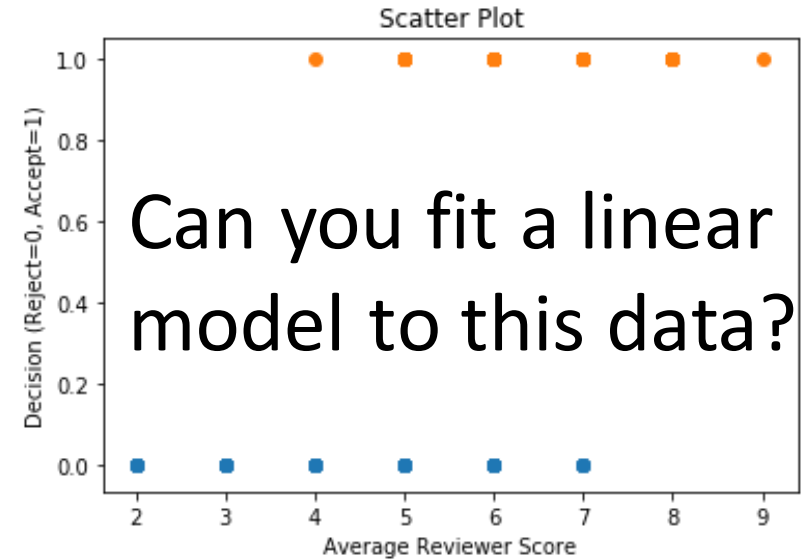
"Supervised Learning (Classification)"

Binary Classification

Binary Classification Task (T):

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

“Categorical
variable”



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

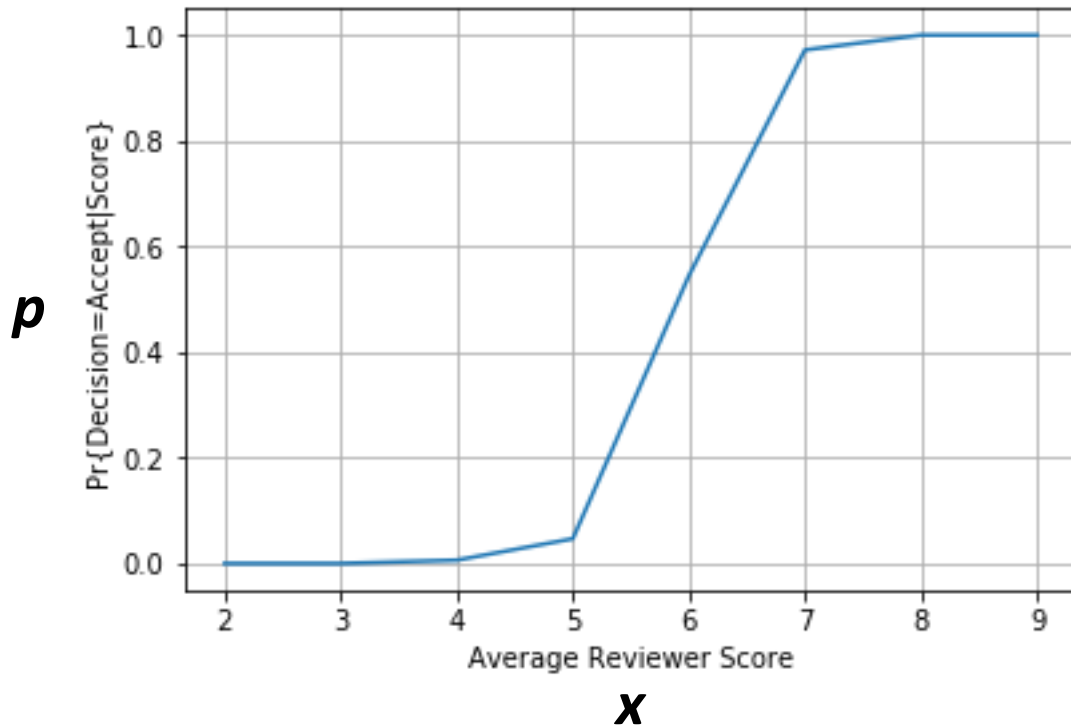
TL;DR	_bibtex	abstract	authorids	authors	conf_1	conf_2	conf_3	decision	review	review_1	review_2	review_3	title
None	@article{\nsharma2018hyperedge2vec,\n\ttitle={H...	Data structured in form of overlapping or non-...	[sharm170@umn.edu, srjoty@ntu.edu.sg, himanshu...	[Ankit Sharma, Shafiq Joty, Himanshu 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,\n\ttitle={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,\n\ttitle={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,\n\ttitle={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):

$\Pr\{\text{Decision}=\text{Accept} \mid \text{Score}\}$

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



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

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

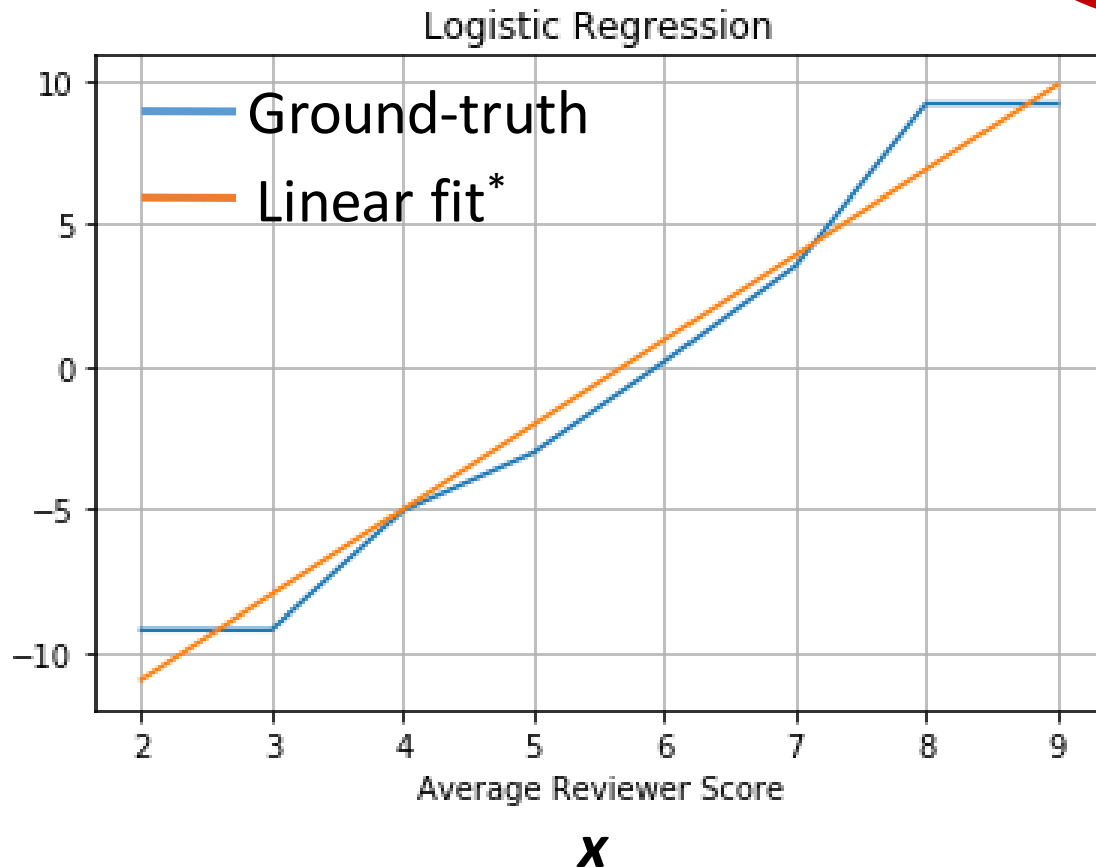
- Is this a good idea?
 - Probability p is always bounded between $[0,1]$

Logistic Regression

Binary Classification Task (T):

“Logits” Function

- Consider the following 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!**

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

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

Logistic Regression

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



$\Pr\{\text{Decision}=\text{Accept} \mid \text{Score}\}$



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

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

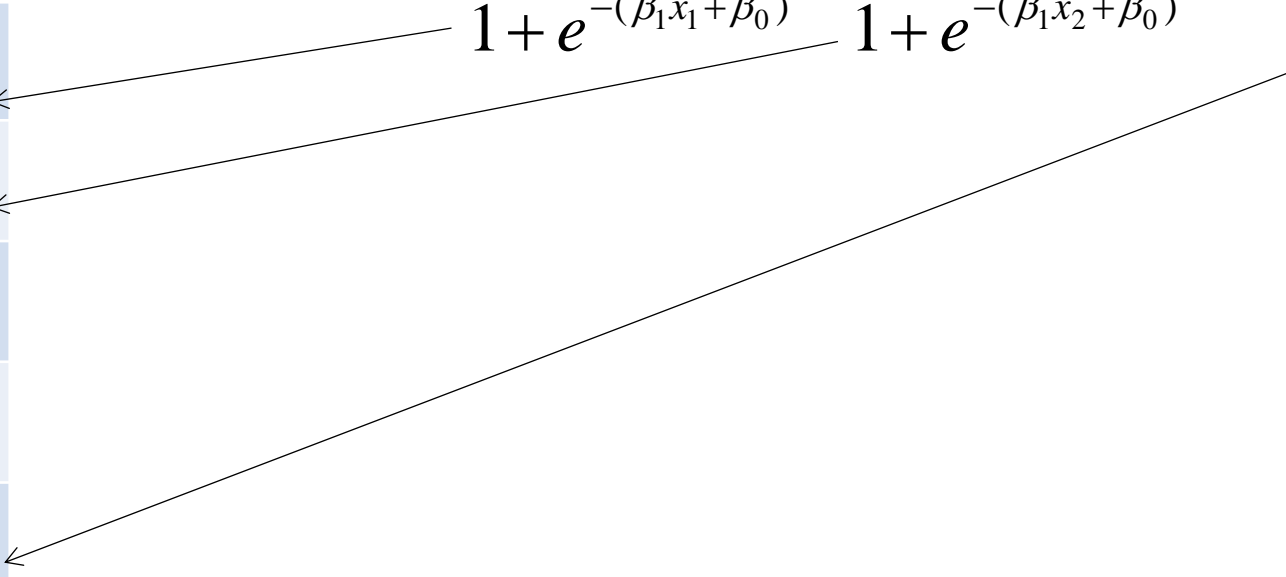
$$1 - p = \frac{e^{-(\beta_1 x + \beta_0)}}{1 + e^{-(\beta_1 x + \beta_0)}}$$

How do we find the model parameters β_1 and β_0 ?

Model Estimation

- We will use an approach referred to as **Maximum Likelihood Estimation** (MLE)
 - Let's assume that the model(i.e., β_1 and β_0) is magically known. Consider the training dataset below. What is the likelihood that the dataset came from our model?

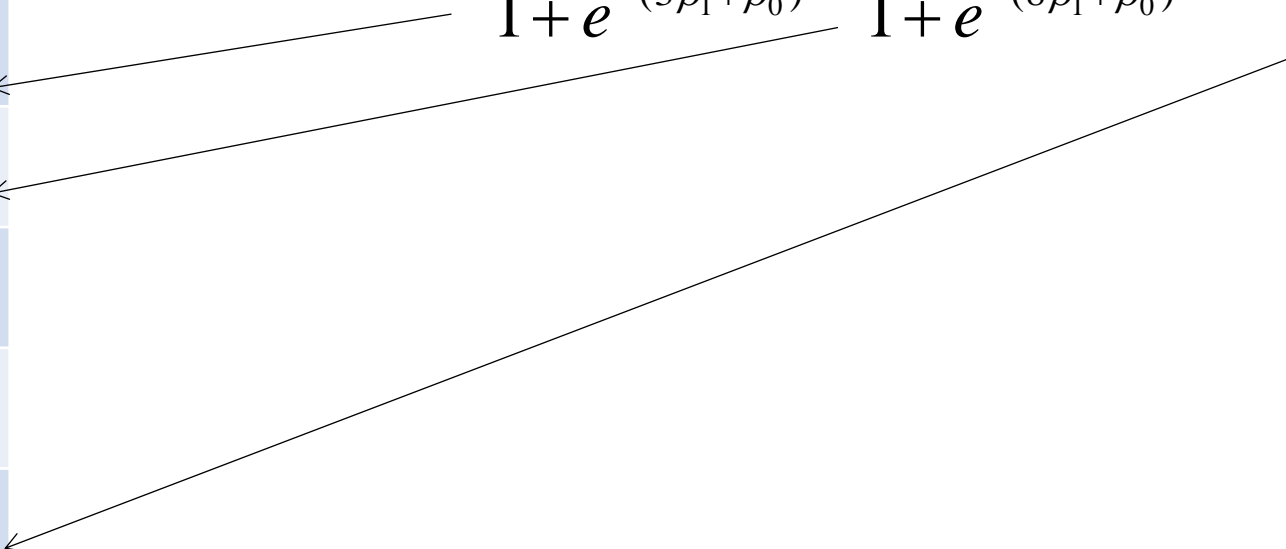
#	X	Y
1	$x_1 = 3$	$y_1 = 0$
2	$x_2 = 8$	$y_2 = 1$
..		
..		
N	$x_N = 6$	$y_N = 1$

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


Model Estimation

- We will use an approach referred to as **Maximum Likelihood Estimation** (MLE)
 - Let's assume that the model(i.e., β_1 and β_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., β_1 and β_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	$x_N = 6$	$y_N = 1$

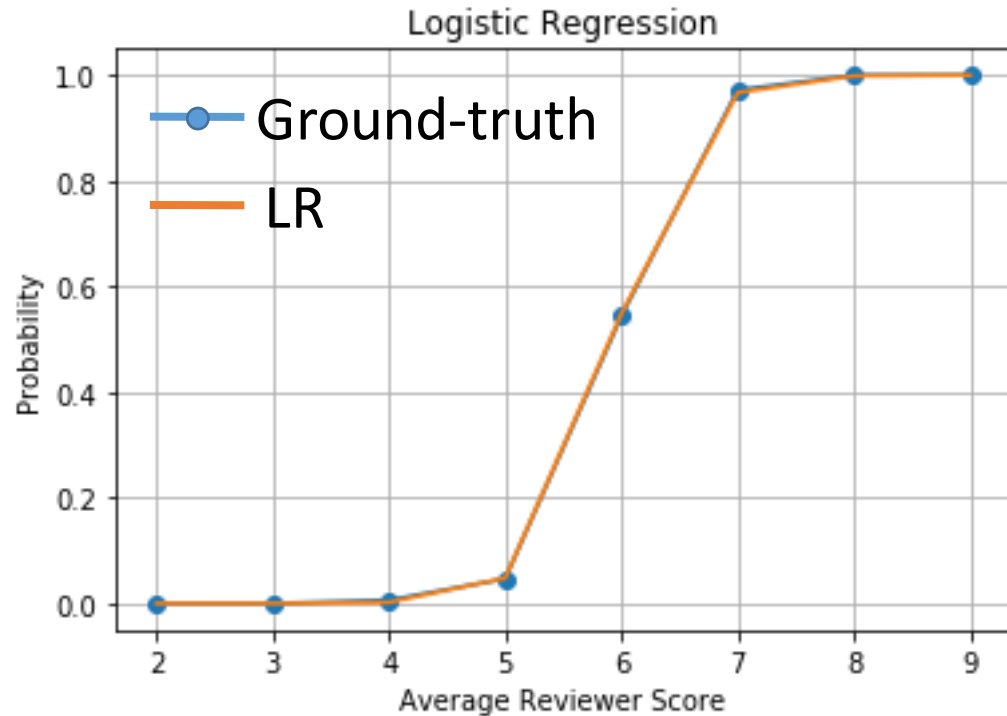
$$\text{Log-Likelihood} = \log\left(\frac{e^{-(3\beta_1 + \beta_0)}}{1 + e^{-(3\beta_1 + \beta_0)}}\right) + \log\left(\frac{1}{1 + e^{-(8\beta_1 + \beta_0)}}\right) + \dots \log\left(\frac{1}{1 + e^{-(6\beta_1 + \beta_0)}}\right)$$

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

Find β_1 and β_0 that maximize g
(or minimize the “loss” $-g$)

$$\text{Loss}(\beta_1, \beta_0) = -g(\beta_1, \beta_0)$$

We Won't Worry About How



```
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.

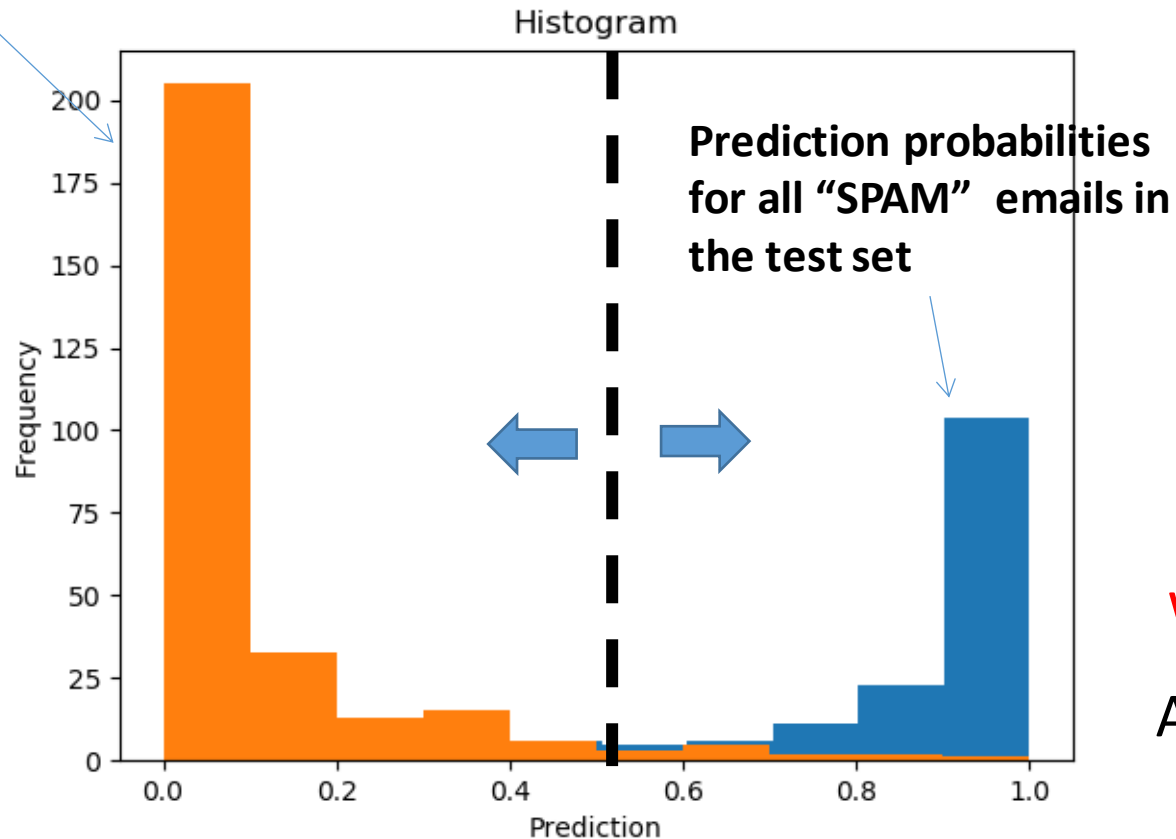
- **57 Real or integer valued features**
- **Binary output class**

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

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%

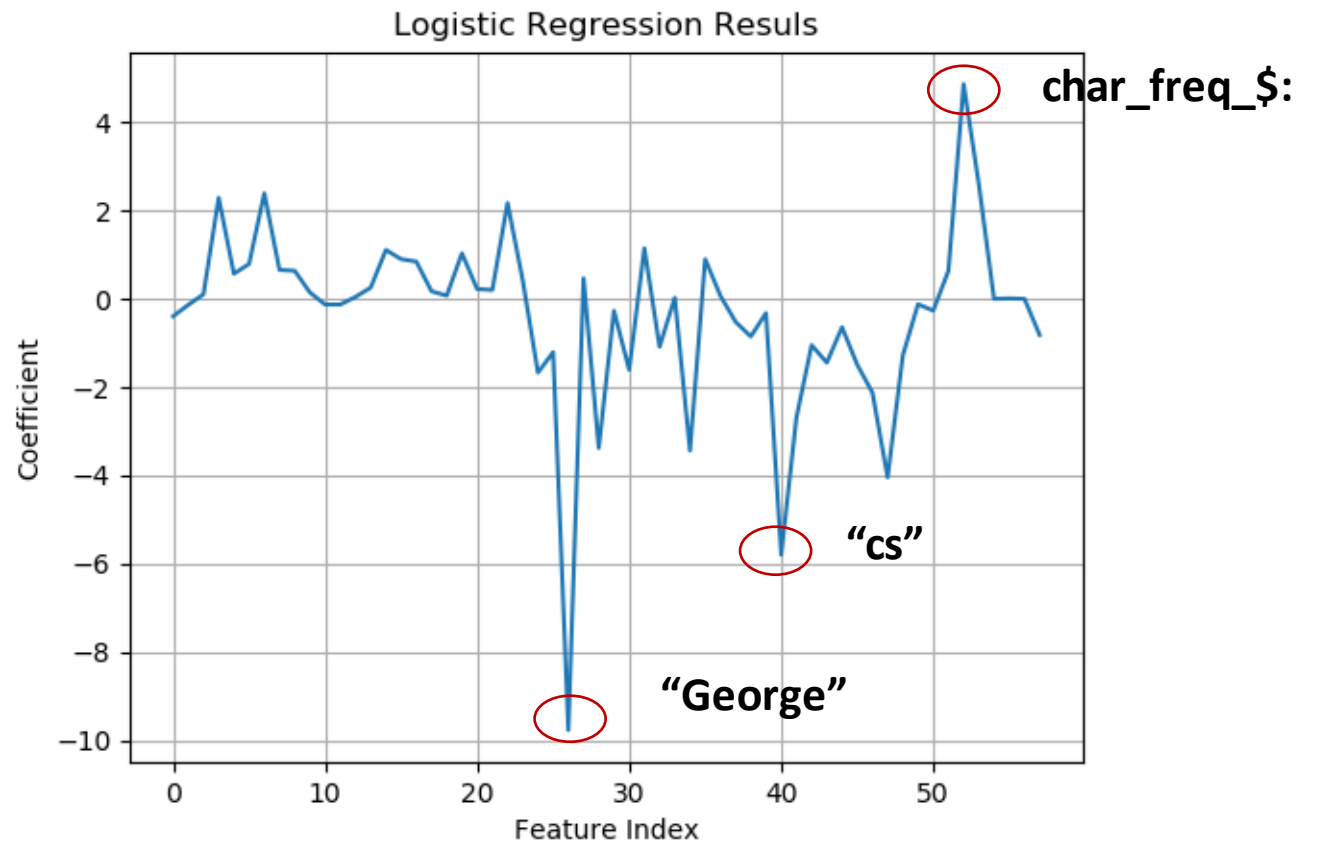
Which Features Matter?

Our Model:

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

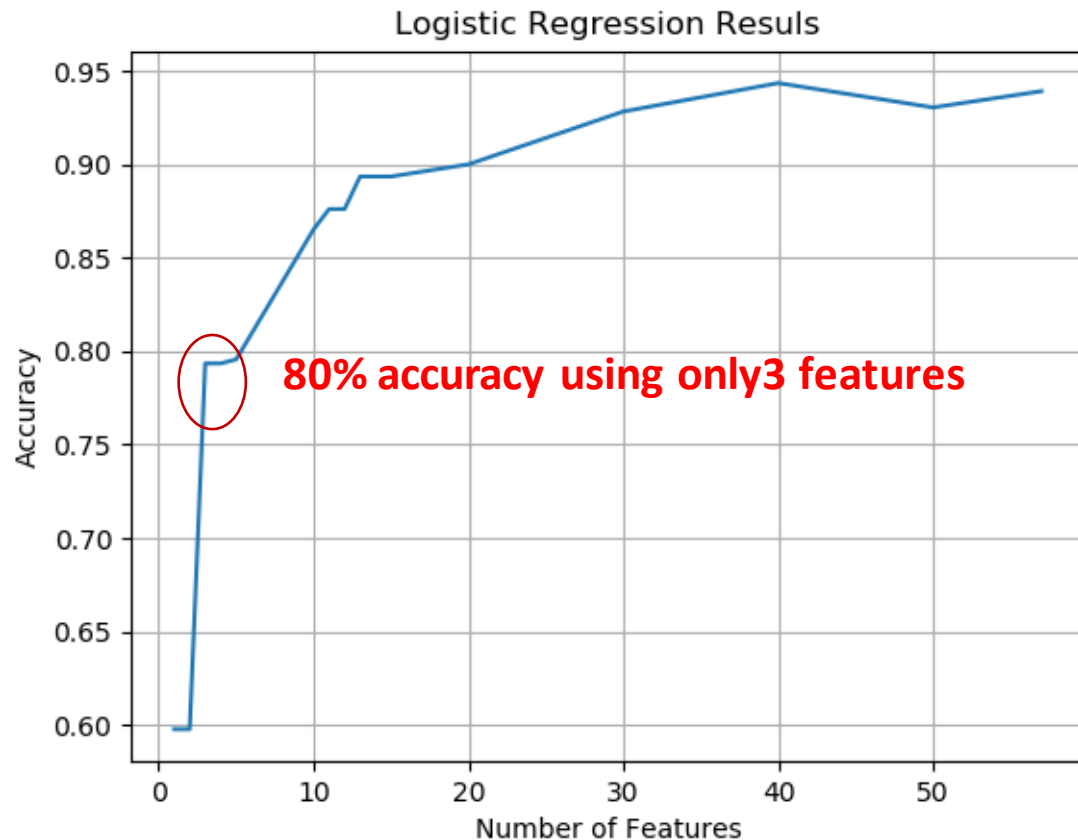
Reasonable hypothesis: features with larger absolute values of β matter more.

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

p	L _p Norm	Interpretation
0	$\ \mathbf{x}\ _0 = (\sum x_i ^0)^{1/0}$	Number of Non-zero Entries
1	$\ \mathbf{x}\ _1 = (\sum x_i)$	Sum of absolute values
2	$\ \mathbf{x}\ _2 = (\sum x_i ^2)^{0.5}$	Root mean square
∞	$\ \mathbf{x}\ _\infty = (\sum x_i ^\infty)^0$	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
(LASSO)**

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

**L2 Regularization
(Ridge)**

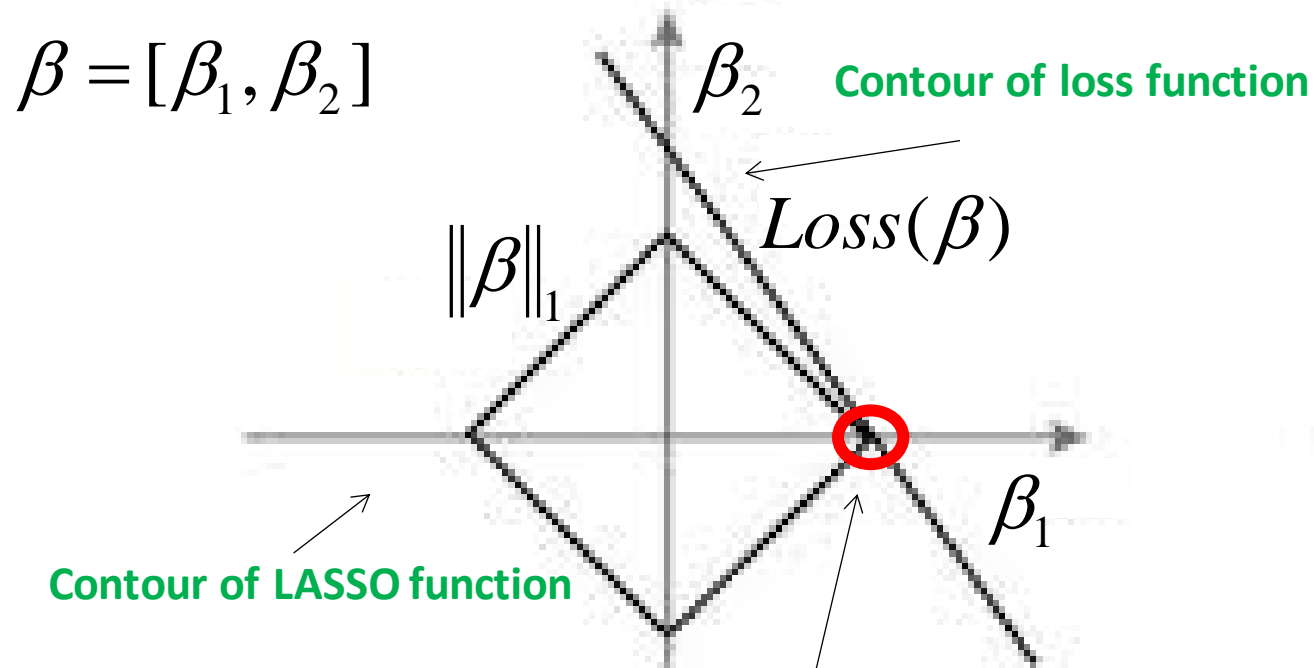
$$\hat{\beta} = \min_{\beta} \{ Loss(\beta) + c \|\beta\|_2 \}$$

We are penalizing
“large” coefficients.

But why?

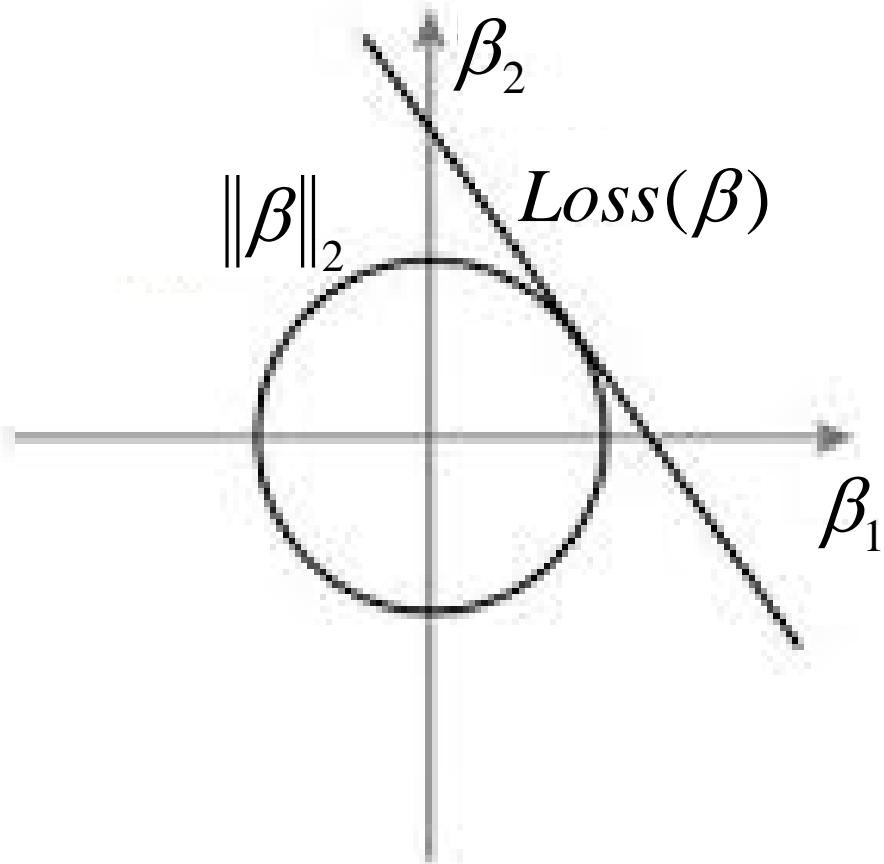
LASSO and Ridge Regularization

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

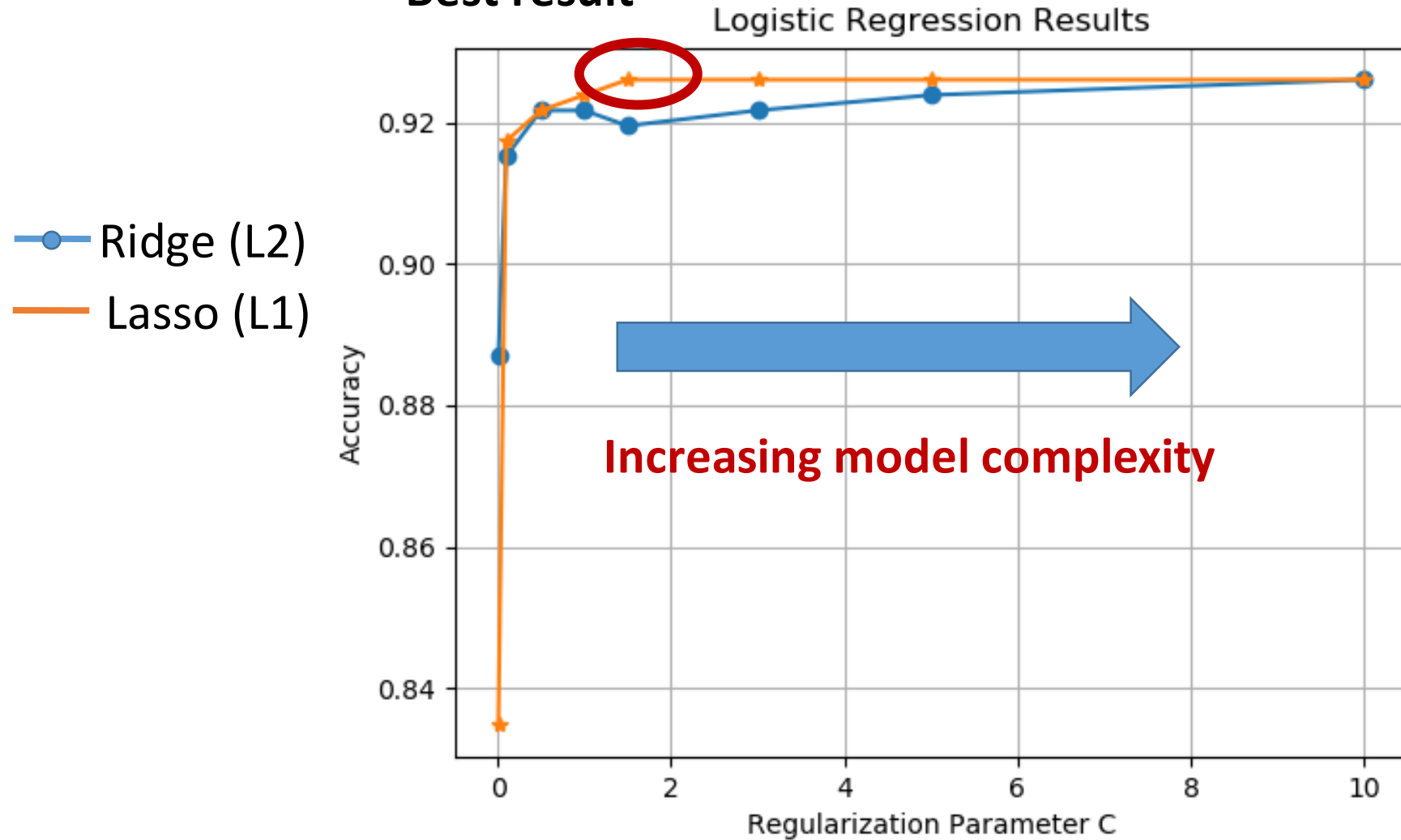
Inverse of regularization strength; must be a positive float. Like in support vector machines, smaller values specify stronger regularization.

Which regularization function to use?

How should we select c ?

Impact of C

Best result



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↓	<i>Class_pos</i>	<i>Class_neg</i>
<i>Class_pos</i>	TP	FN
<i>Class_neg</i>	FP	TN

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

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

Hard Decisions

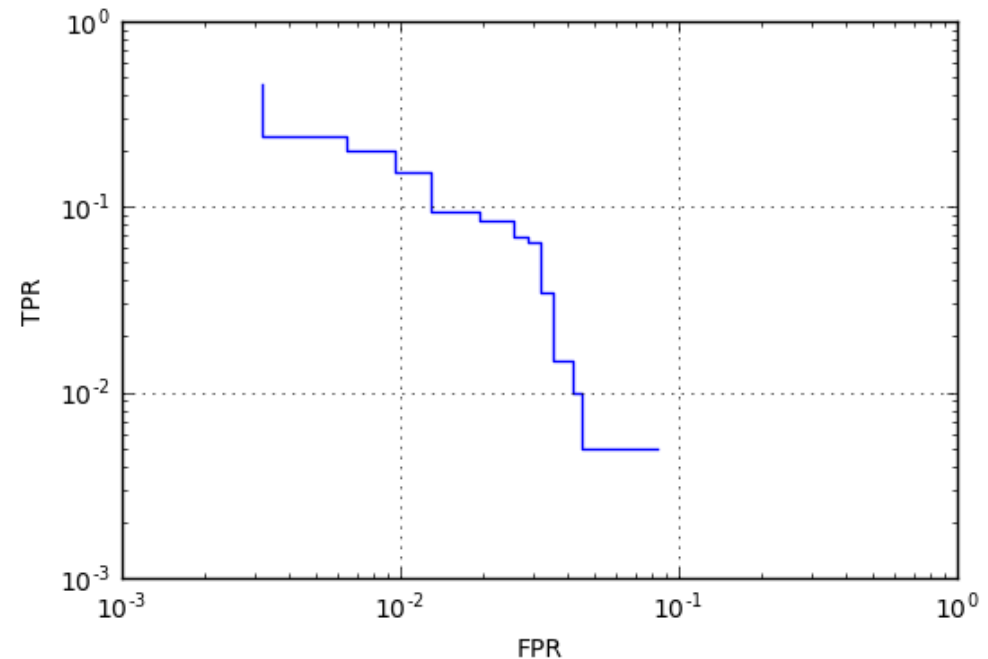
- Logistic classifier outputs a **soft** label: $P(y = 1|x) \in [0,1]$
 - $P(y = 1|x) \approx 1 \Rightarrow y = 1$ more likely
 - $P(y = 0|x) \approx 1 \Rightarrow y = 0$ more likely
- Can obtain a **hard label** by **thresholding**:
 - Set $\hat{y} = 1 \Leftrightarrow P(y = 1|x) > t$
 - t = Threshold
- How to set threshold?
 - Set $t = \frac{1}{2} \Rightarrow$ Minimizes overall error rate
 - Increasing $t \Rightarrow$ Decreases false positives
 - Decreasing $t \Rightarrow$ Decreases missed detections

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



Linear Classifiers

- LR based spam detector that we discussed is an example of a **Linear Classifier**

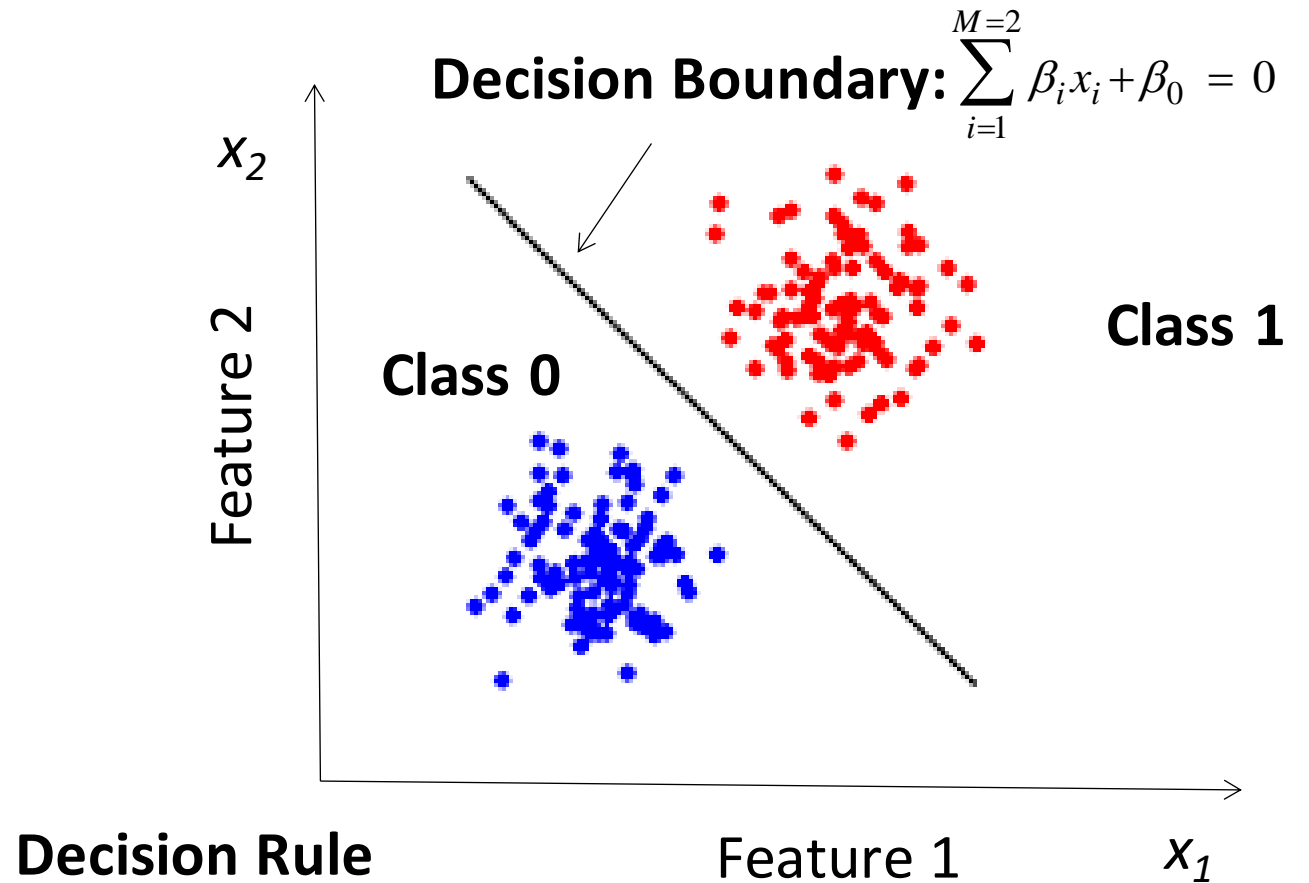
Recall that:

$$p_{Class1} = \frac{1}{1 + e^{-\left(\sum_{i=1}^M \beta_i x_i + \beta_0\right)}}$$

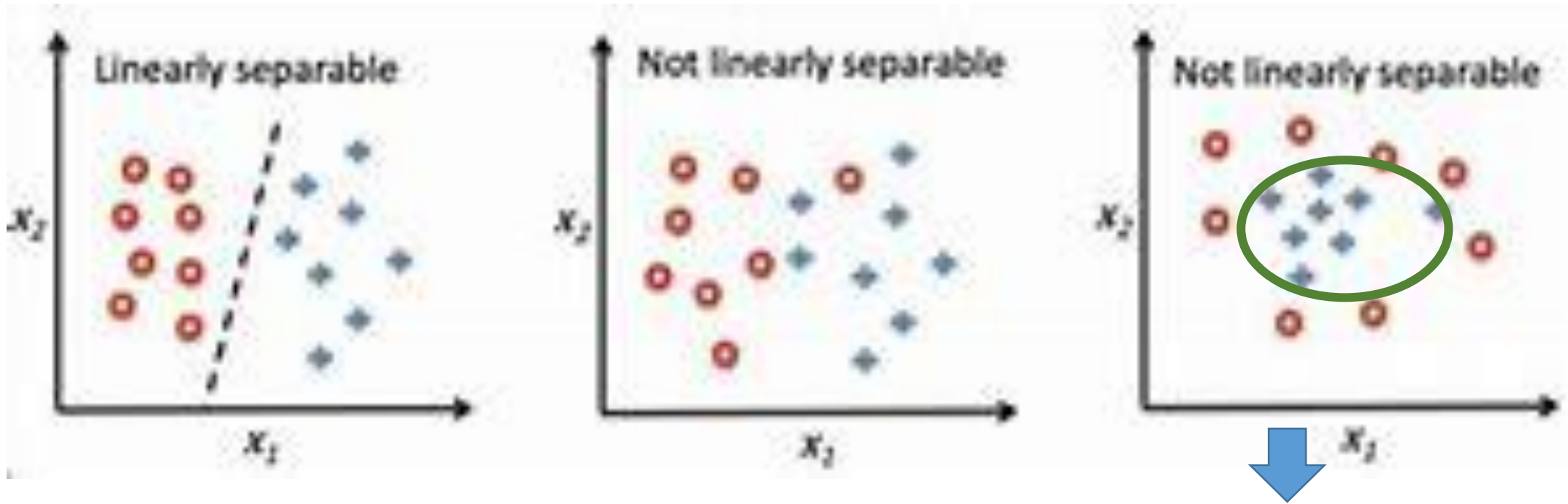
$$x \in \text{Class 1} \text{ iff } p_{Class1} \geq 0.5$$

$$\Rightarrow e^{-\left(\sum_{i=1}^M \beta_i x_i + \beta_0\right)} \leq 1$$

$$\Rightarrow \left(\sum_{i=1}^M \beta_i x_i + \beta_0\right) \geq 0$$



Linear Separability

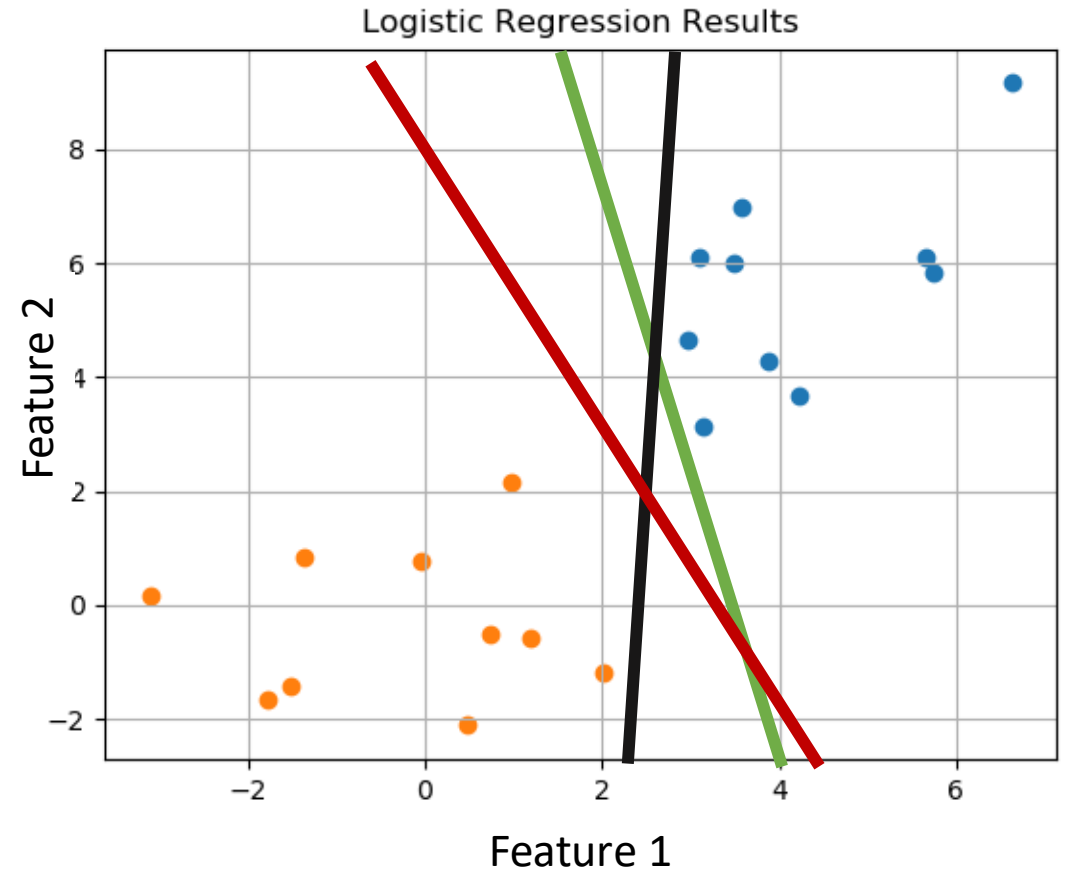


Is there a transformation of the features for which this data is linearly separable?

Support Vector Machines (SVMs)

- Recall that LR attempts to minimize the log-likelihood of the training data under the LR model
 - LR has a probabilistic interpretation
- SVMs on the other hand have a **geometric interpretation**
 - Goal: maximize margin to closest training data sample

All 3 classifiers have 100% training accuracy



But which of the 3 classifiers would you pick?

SVMs: Maximize Margin

- Assume data is linearly separable
- Goal: maximize margin to closest training data sample

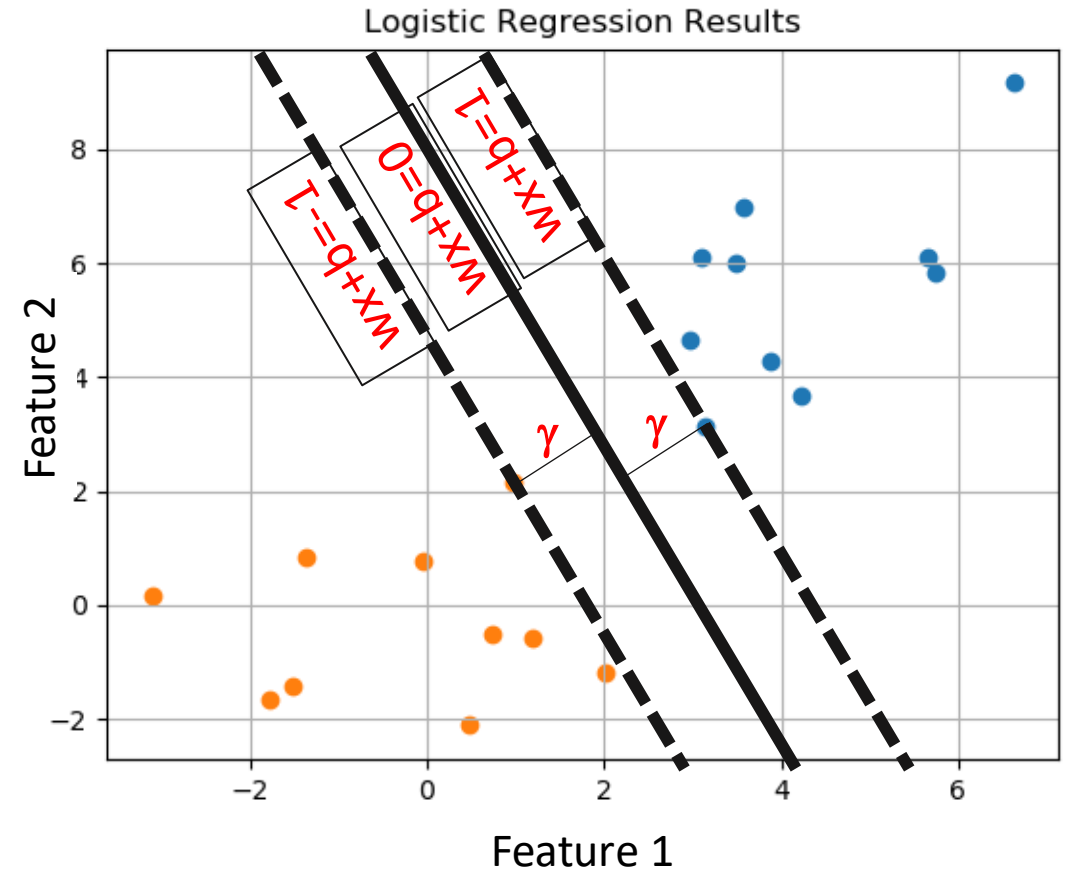
- Let the closest points on either side lie on :

$$wx + b = \pm 1$$

- What is the value of the margin γ ?

$$\lambda = \frac{1}{\|w\|_2}$$

Maximize



SVMs: Optimization

- Optimization problem

$$\max \frac{1}{\|w\|_2} \Rightarrow \min \|w\|_2$$

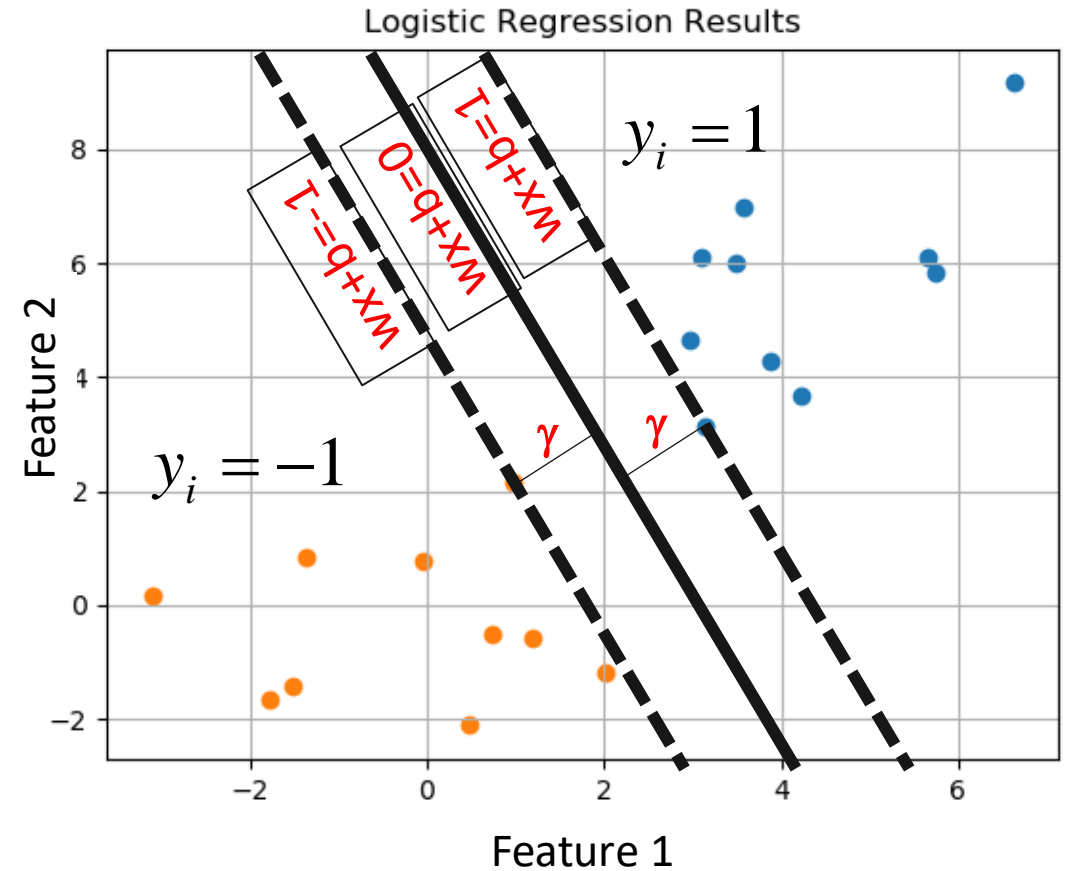
$$wx_i + b \geq 1 \quad y_i = 1$$

$$wx_i + b \leq -1 \quad y_i = -1$$

- Or equivalently

$$\min \|w\|_2$$

$$y_i(wx_i + b) \geq 1$$



SVMs: Non-linearly Separable

- Make the constraints “soft” by introducing an “margin” ε_i

$$\text{Min } \sum_i \varepsilon_i + \alpha \|w\|_2$$

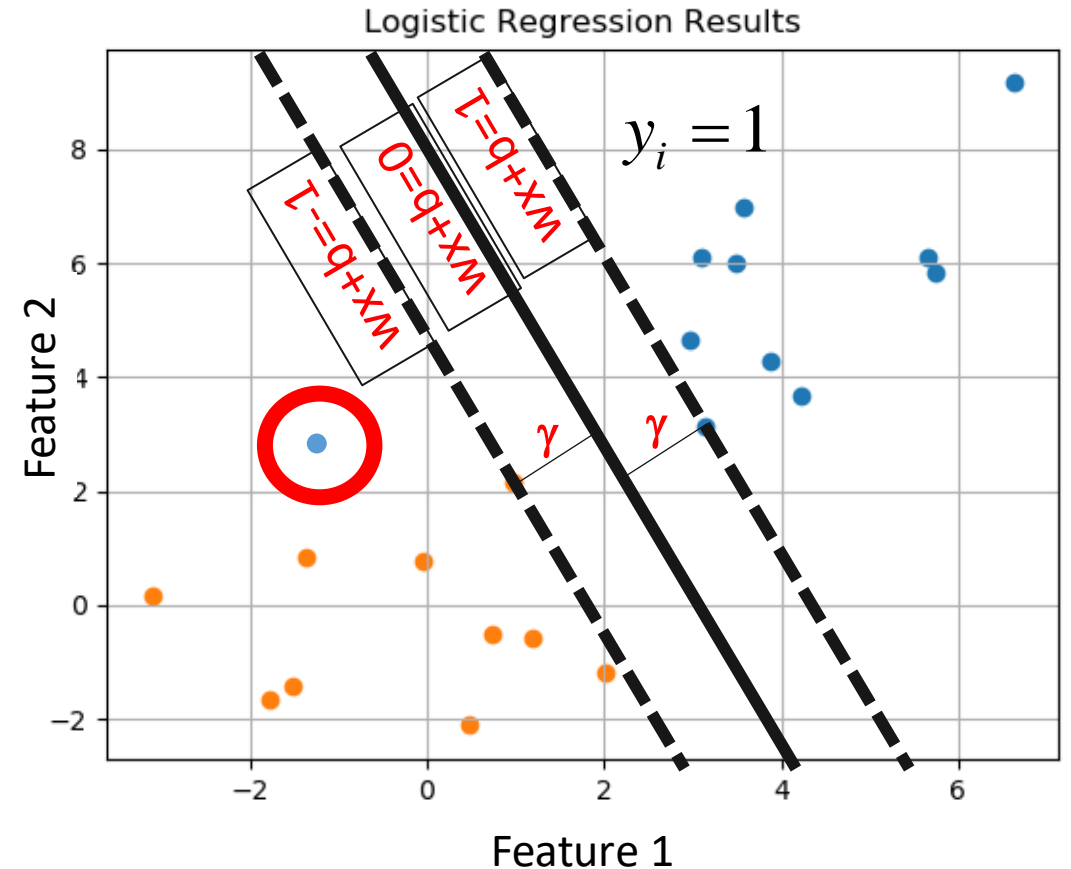
$$y_i(wx_i + b) \geq 1 - \varepsilon_i$$

$$\varepsilon_i \geq 0$$

Points can be misclassified with “cost” ε_i

- Note that in any solution to this problem

$$\varepsilon_i = \max\{0, 1 - y_i(wx_i + b)\}$$



SVMs: Putting it All Together

- SVMs minimize the following objective

Regularization constant α

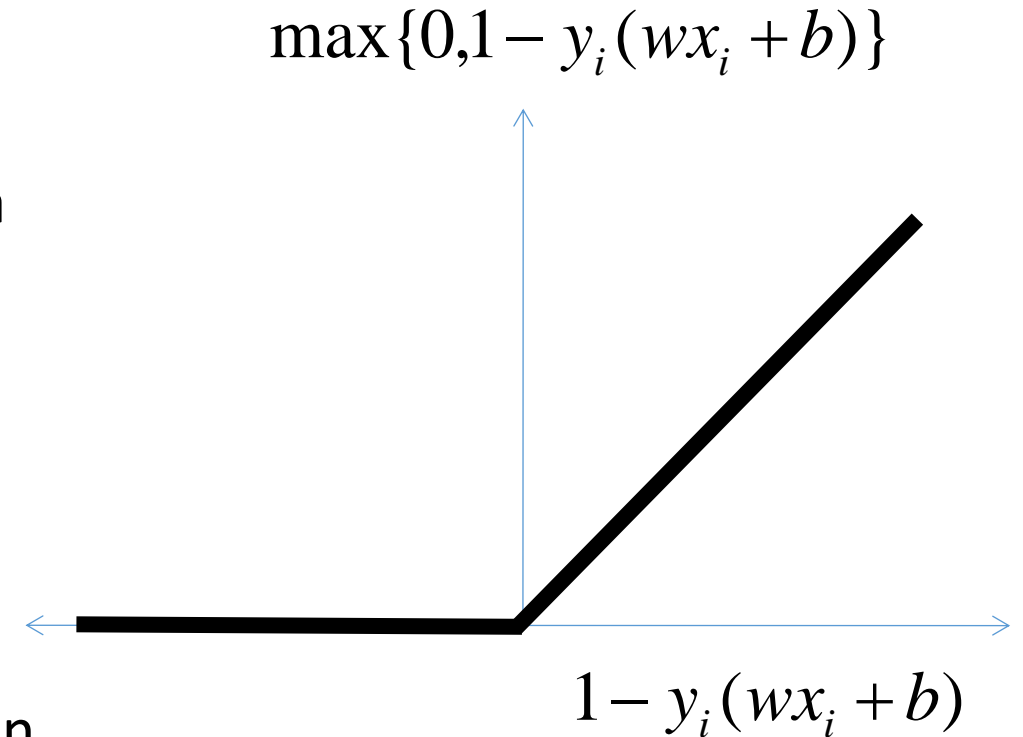
Increasing α results in greater margins at the expense of misclassifications on training data

$$\max\{0, 1 - y_i(wx_i + b)\} + \alpha \|w\|_2$$

“Hinge Loss”

Seeks to minimize misclassifications

Tries to maximize margin
(L_2 regularizer!)



Relation to Logistic Regression

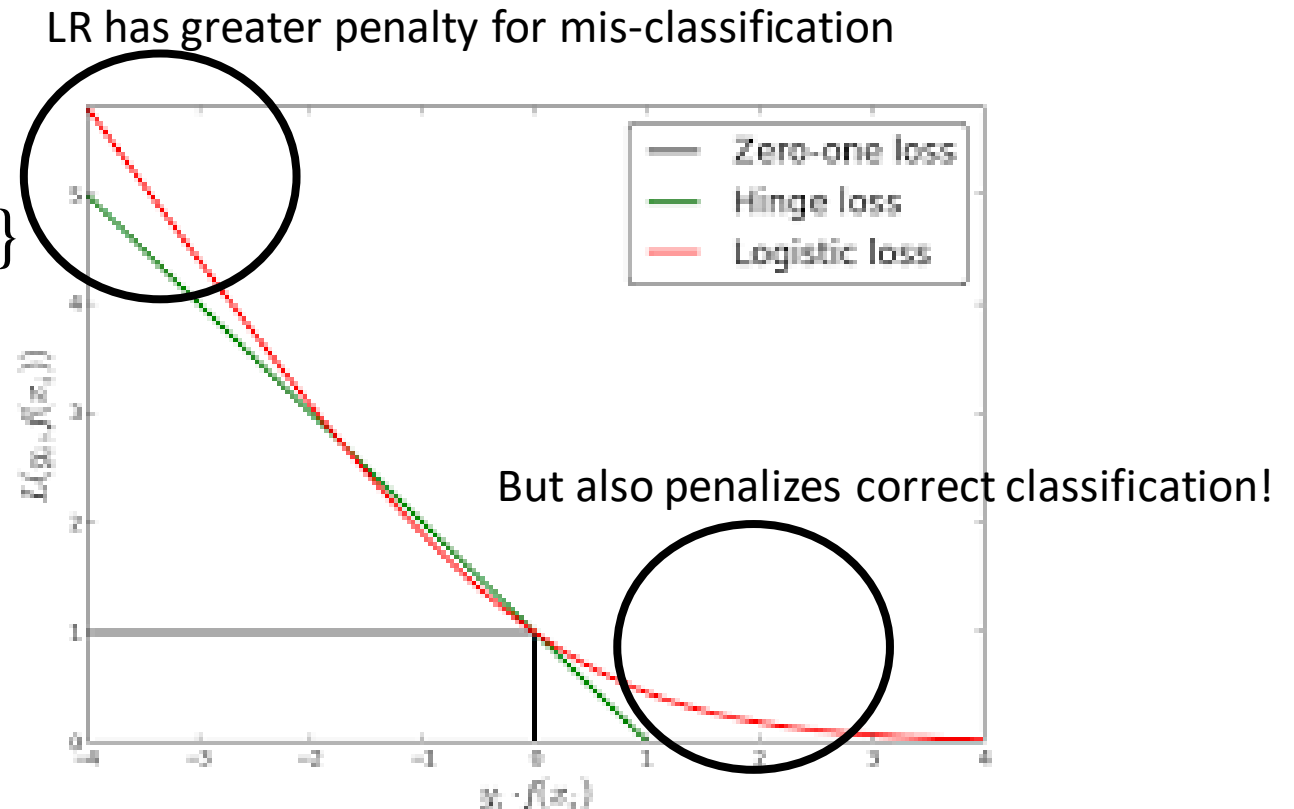
- LR and SVM differ only on the loss function they use
 - “Log-loss” versus “hinge-loss”

Hinge Loss:

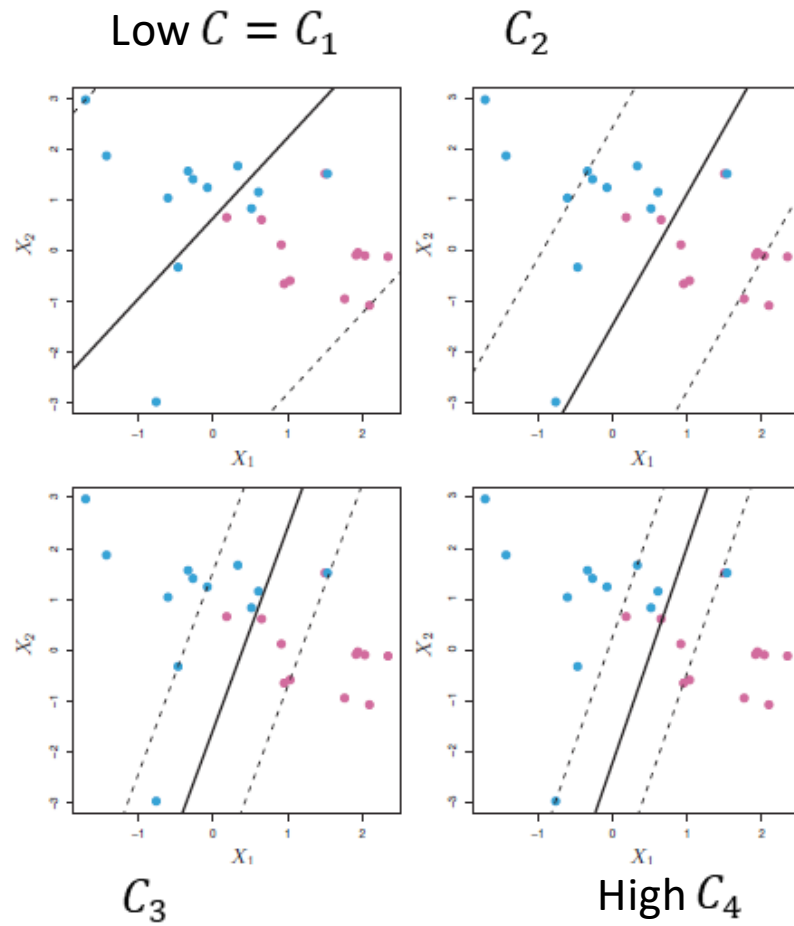
$$\max\{0, 1 - y_i(wx_i + b)\} = \max\{0, 1 - q\}$$

Log Loss:

$$-\log_2\left(\frac{1}{1 + e^{-y_i(wx_i + b)}}\right) = \log_2(1 + e^{-q})$$



Illustrating Effect of C

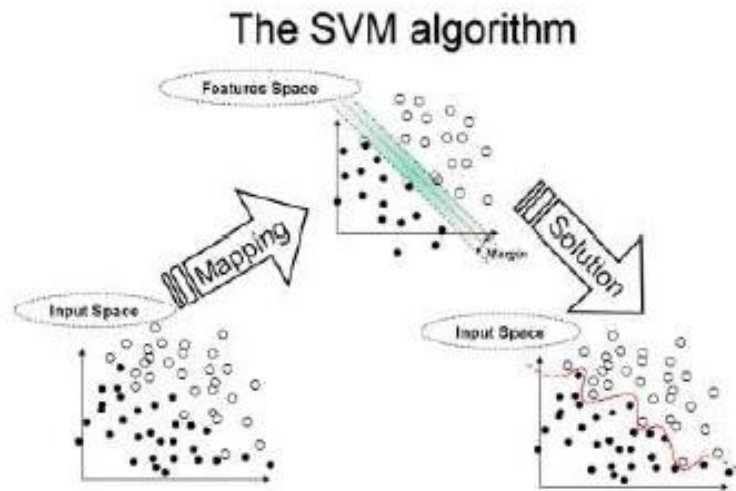


- Fig. 9.7 of ISL
 - Note: C has opposite meaning in ISL than python
 - Here, we use python meaning
- Low C :
 - Many violations of margin.
 - Many more SVs
 - Reduces variance by using more samples

Transform Problem

- Transform problem: replace x with $\phi(x)$
 - Enables more rich classifiers
 - Examples: polynomial classification $\phi(x) = [1, x, x^2, \dots, x^{d-1}]$
- Tries to find separation in a **feature** space

From <https://www.dtreg.com/solution/view/20>



Kernel Trick

- Classifier is:

- $z = b + \mathbf{w}^T \mathbf{x} = b + \sum_{i=1}^N \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}), \quad K(\mathbf{x}_i, \mathbf{x}) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x})$

- $\hat{y} = \text{sign}(z) = \begin{cases} 1 & \text{if } z > 0 \\ -1 & \text{if } z < 0 \end{cases}$

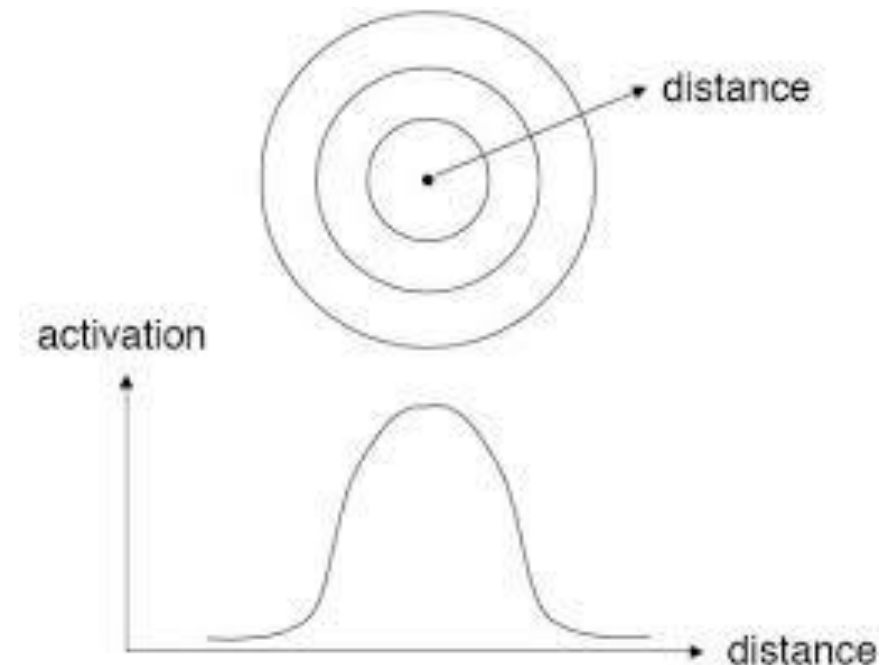
- Do not need to explicitly compute $\phi(\mathbf{x})$

- Can directly compute kernel $K(\mathbf{x}_i, \mathbf{x})$

- Provided kernel corresponds to some $\phi(\mathbf{x})$

Understanding the Kernel

- Kernel function $K(\mathbf{x}_i, \mathbf{x})$:
 - measures “distance” between new sample \mathbf{x} and training data \mathbf{x}_i
 - $K(\mathbf{x}_i, \mathbf{x})$ large $\Rightarrow \mathbf{x}_i, \mathbf{x}$ close
 - $K(\mathbf{x}_i, \mathbf{x}) \approx 0 \Rightarrow \mathbf{x}_i, \mathbf{x}$ far
- Linear discriminant $z = b + \sum_{i=1}^N \alpha_i y_i K(\mathbf{x}_i, \mathbf{x})$
 - Weighs sample \mathbf{x}_i that are close to \mathbf{x}

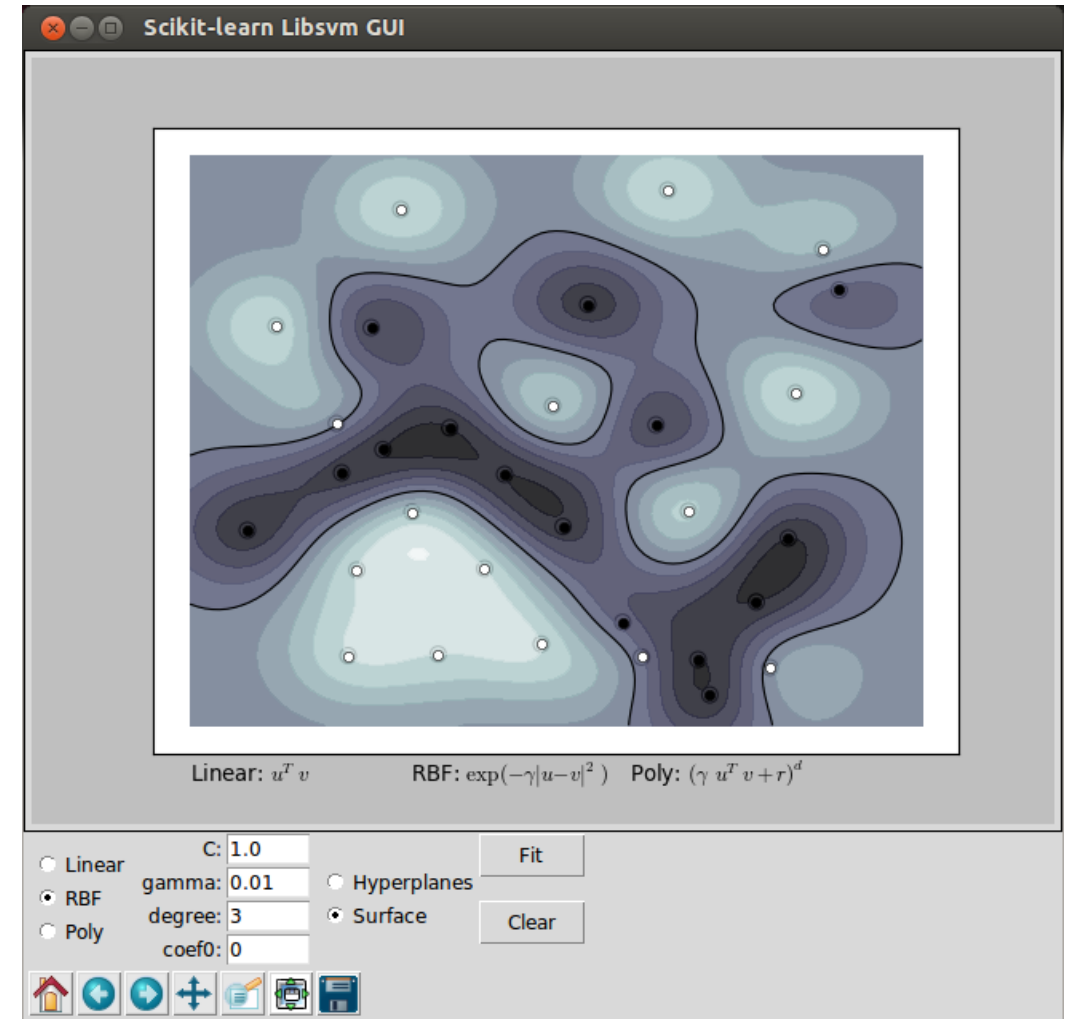


Possible Kernels

- Radial basis function:

$$K(x, x') = \exp[-\gamma \|x - x'\|^2]$$

- $1/\gamma$ indicates width of kernel
- Polynomial kernel: $K(x, x') = |x^T x|^d$
 - Typically $d=2$



Spam Detection: Features

- Recall features used in the UCI Spam database

48 continuous real [0,100] attributes of type word_freq_WORD

- Even easier way to encode features:
 - $x_i = 1$ if term i appears in a document; 0 otherwise
 - Boolean features
- Assume M Boolean features, $x = (x_1, x_2, \dots, x_M)$
 - We want to map this M -dimensional Boolean input to a Boolean output y
 - *Thoughts?*
 - Instead of using LR or SVM we will start with an even simpler approach referred to as “Naïve Bayes”

Ref: Metsis, Vangelis, Ion Androutsopoulos, and Georgios Paliouras. "Spam filtering with naive bayes-which naive bayes?." In *CEAS*, vol. 17, pp. 28-69. 2006.

Naïve Bayes for Spam Filtering

- Assume $M=1$ Boolean feature, $x = (x_1, x_2, \dots, x_M)$
- Each email is either $\{s=\text{spam}, l=\text{legit}\}$
- We begin by computing:

$$P\{spam \mid x\} = \frac{P\{x \mid spam\} * P\{spam\}}{P\{x\}}$$

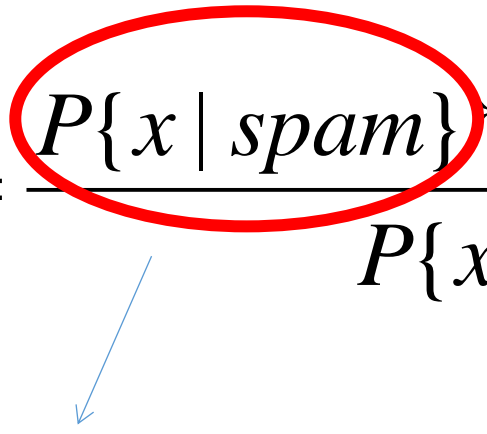
Bayes Rule

$$P\{A \cap B\} = P\{A \mid B\} * P\{B\}$$

Ref: Metsis, Vangelis, Ion Androutsopoulos, and Georgios Paliouras. "Spam filtering with naive bayes-which naive bayes?." In *CEAS*, vol. 17, pp. 28-69. 2006.

Naiive Bayes for Spam Filtering

- We begin by computing:

$$P\{spam | x\} = \frac{P\{x | spam\} * P\{spam\}}{P\{x\}}$$


$$P\{x_1, x_2, \dots, x_M | spam\} = P\{x_1 | spam\} * P\{x_2 | spam\} * \dots * P\{x_M | spam\}$$

Assuming that term occurrences are independent (given class)!

Is this a reasonable assumption?

Naïve Bayes for Spam Filtering

$$P\{x_1 | spam\} * P\{x_2 | spam\} * .. * P\{x_M | spam\}$$



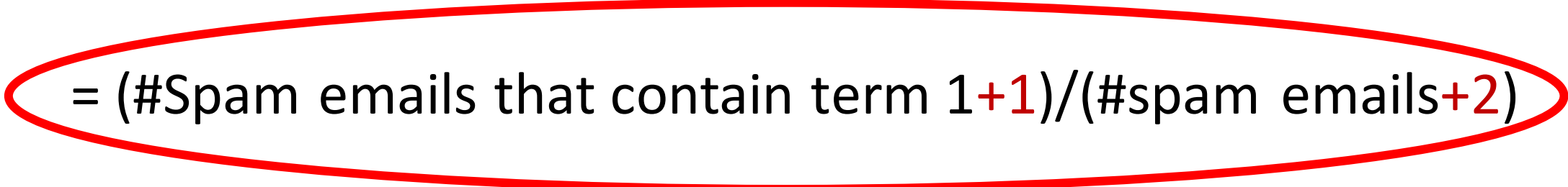
How do we estimate this from the training dataset?

$$P\{x_1 = 1 | spam\} = (\text{\#Spam emails that contain term 1}) / (\text{\#spam emails})$$



What happens if term 1 never occurred in any spam email in the training dataset?

Laplacian Smoothing

$$P\{x_1 = 1 \mid \textit{spam}\} = \cancel{(\# \text{Spam emails that contain term 1}) / (\# \text{spam emails})}$$

$$= (\# \text{Spam emails that contain term } 1+1) / (\# \text{spam emails}+2)$$

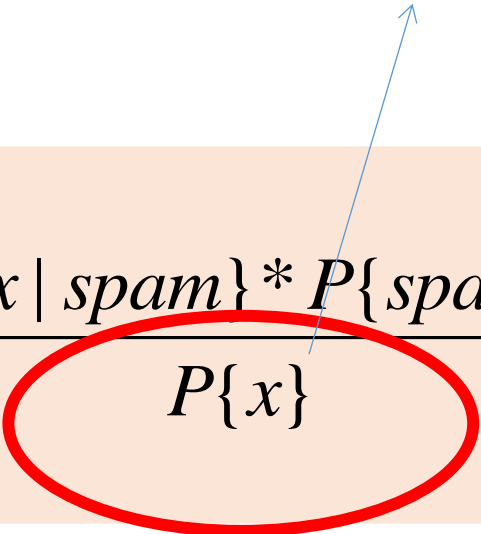
Equivalent to assuming two additional spam emails in the training dataset, of which one contains all terms and the other is empty

$$P\{x_1 = 0 \mid \textit{spam}\} = 1 - P\{x_1 = 1 \mid \textit{spam}\}$$

$$= (\# \text{Spam emails that don't contain term } 1+1) / (\# \text{spam emails}+2)$$

Naiive Bayes for Spam Filtering

$$P\{x\} = P\{spam\} * P\{x | spam\} + P\{legit\} * P\{x | legit\}$$


$$P\{spam | x\} = \frac{P\{x | spam\} * P\{spam\}}{P\{x\}} \quad \text{Vs.} \quad P\{legit | x\} = \frac{P\{x | legit\} * P\{legit\}}{P\{x\}}$$

Or:

$$P\{spam | x\} \geq threshold$$

References

- Cormack, Gordon V. "Email Spam Filtering: A Systematic Review." *Foundations and Trends® in Information Retrieval* 1.4 (2008): 335-455. https://www.ccs.neu.edu/home/vip/teach/IRcourse/IR_surveys/spam-filtering.pdf
- Metsis, Vangelis, Ion Androutsopoulos, and Georgios Paliouras. "Spam filtering with naive bayes-which naive bayes?." In *CEAS*, vol. 17, pp. 28-69. 2006.