# Logistic Regression for Text Classification

CS 485, Fall 2024
Applications of Natural Language Processing

Brendan O'Connor
College of Information and Computer Sciences
University of Massachusetts Amherst

*[With slides from Ari Kobren and SLP3]*

# BOW linear model for text classif.

- Problem: classify doc $d$ into one of $k \in 1..K$ classes

- Parameters: For each class $k$, and word type $w$, there is a *word weight*

$$\beta_{w,k}$$

$$\beta \in \mathbb{R}^{|V| \times K}$$

|  | $k=NEG$ | $k=POS$ |
|---|---|---|
| dog | $-0.1$ | $0.1$ |
| happy | $-1.4$ | $3.1$ |
| ? | | |

- Representation: bag-of-words vector of doc $d$'s word counts

$$x = (1, 4, 1, 0, 0 \ldots -)$$

$$x_w \in \mathbb{R}^V$$

- Prediction rule: choose class $y$ with highest score

$$Score(y) = \sum_{w \in V} x_w \beta_{w,k}$$

$$\hat{y} = \underset{k \in 1..K}{\text{argmax}} \; Score(k)$$

# Keyword count as linear model

- Problem: classify doc $d$ into one of $k \in 1..K$ classes

(of the <u>happy</u> mare was <u>grief</u> :-) )

- Parameters: For each class $k$, and word type $w$, there is a *word weight*

$$L_{POS} = \{happy, great \cdots\}$$
$$L_{NEG} = \{\cdots\cdots\cdots\}$$
$$L_k = \{\cdots\cdots\cdots\}$$
$$\beta_{w,k} = \mathbb{1}\{w \in L_k\}$$

- Representation: bag-of-words vector of doc $d$'s word counts

$$x_w$$

- Prediction rule: choose class $y$ with highest score

$$\implies \sum_{w \in V} x_w \beta_{w,k} = \text{how many words from doc.}$$
$$\text{are in lexicon "k"}$$

$$\implies \underset{k}{argmax}\ score(k) = \text{highest lex. count}$$
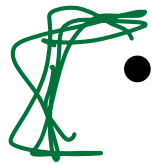
# Naive Bayes as linear model

- Problem: classify doc $d$ into one of $k \in 1..K$ classes

- Parameters: For each class $k$, and word type $w$, there is a *word weight*

- Representation: bag-of-words vector of doc $d$'s word counts

- Prediction rule: choose class $y$ with highest score

# Linear classification models

- The foundational model for machine learning-based NLP!

- Examples
  - The humble "keyword count" classifier (no ML)
  - Naive Bayes ("generative" ML)

- Today: **Logistic Regression**
  - a linear classification model, trained to be good at *prediction*
  - allows for *features*
  - used within more complex models (neural networks)

# Motivation: feature engineering

- For Naive Bayes, we used counts of each word in the vocabulary (BOW representation).  But why not also use....
  - Number of words from "CS485 Crowdsource Positive Lexicon"
  - ...from "CS485 Crowdsource Negative Lexicon" ... or another....
  - Phrases?
  - Words/phrases with negation markers?
  - Number of "!" occurrences?
  - or...?

- NB tends to work poorly when there are many potentially repetitive features (why?) → Cond. Indep. assump. is wrong!

# Features! Features! Features!

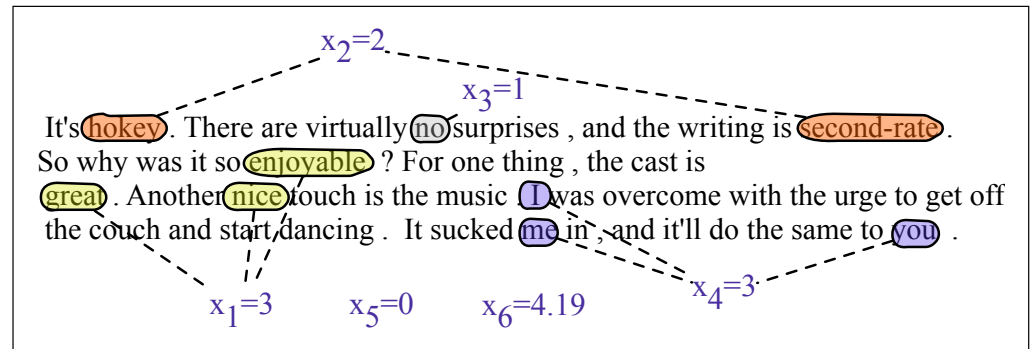| Var | Definition | Value in Fig. 5.2 |
|---|---|---|
| $x_1$ | count(positive lexicon words $\in$ doc) | 3 |
| $x_2$ | count(negative lexicon words $\in$ doc) | 2 |
| $x_3$ | $\begin{cases} 1 & \text{if "no"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$ | 1 |
| $x_4$ | count(1st and 2nd pronouns $\in$ doc) | 3 |
| $x_5$ | $\begin{cases} 1 & \text{if "!"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$ | 0 |
| $x_6$ | ln(word count of doc) | $\ln(66) = 4.19$ |

It's hokey. There are virtually no surprises, and the writing is second-rate. So why was it so enjoyable? For one thing, the cast is great. Another nice touch is the music I was overcome with the urge to get off the couch and start dancing. It sucked me in, and it'll do the same to you.

$x_2=2$   $x_3=1$   $x_1=3$   $x_5=0$   $x_6=4.19$   $x_4=3$

**Figure 5.2**   A sample mini test document showing the extracted features in the vector x.

modate *any arbitrary features*

spend a lot of trying and testing

! This is a place to put linguistics in,

your data.

# Negation

Das, Sanjiv and Mike Chen. 2001. Yahoo! for Amazon: Extracting market sentiment from stock message boards. In Proceedings of the Asia Pacific Finance Association Annual Conference (APFA).
Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. 2002. Thumbs up? Sentiment Classification using Machine Learning Techniques. EMNLP-2002, 79—86.

Add NOT_ to every word between negation and following punctuation:

didn't like this movie , but I



didn't NOT_like NOT_this NOT_movie but I

*[Slide: SLP3]*

# Classification: LogReg (I)

First, we'll discuss **how LogReg works.**

① Binary LogReg

③ Learning

② Multiclass LR

Then, **why** it's set up the way that it is.

Application: **spam filtering**

# Classification: LogReg (I)

- compute **features** (xs)

$$x_i = (\text{count "nigerian", count "prince", count "nigerian prince"})$$

- given **weights** (betas)

$$\beta = (-1.0, \quad -1.0, \quad 4.0)$$

# Classification: LogReg (II)

- Compute the **dot product**

$$z = x^T \beta = x' \beta = \sum_{j=1}^{3} x_j \beta_j$$

$$z = \sum_{i=0}^{} \beta_i x_i = z$$

$$\boxed{z} = x^T \beta = x' \beta$$

$$|x|$$

- Compute the **logistic function** for the label probability

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

$$P(y=1 \mid x) = \frac{e^z}{1 + e^z} = \frac{1}{1 + e^{-z}}$$

$$P(y=0 \mid x) = 1 - P(y=1 \mid x)$$

# LogReg Exercise

features: $\left(\text{count "nigerian", count "prince", count "nigerian prince"}\right)$

$$x = (1, \quad 1, \quad 1)$$

$$\beta = (-1.0, \quad -1.0, \quad 4.0)$$

$$z = x^T\beta = x_1\beta_1 + x_2\beta_2 + x_3\beta_3 = 2$$

$$P(x) = ???$$

$$P(y=1 \mid x) = \frac{1}{1+e^{-2}} = \frac{1}{1+e^{-2}} = 0.88$$

# Classification: Dot Product

$$z = \sum_{j=0}^{\text{Nfeat}} \beta_j x_{ij}$$

Nfeat = |V| for BOW

Weighted sum of feature values

"Linear Model"

# Why the **logistic function**?

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

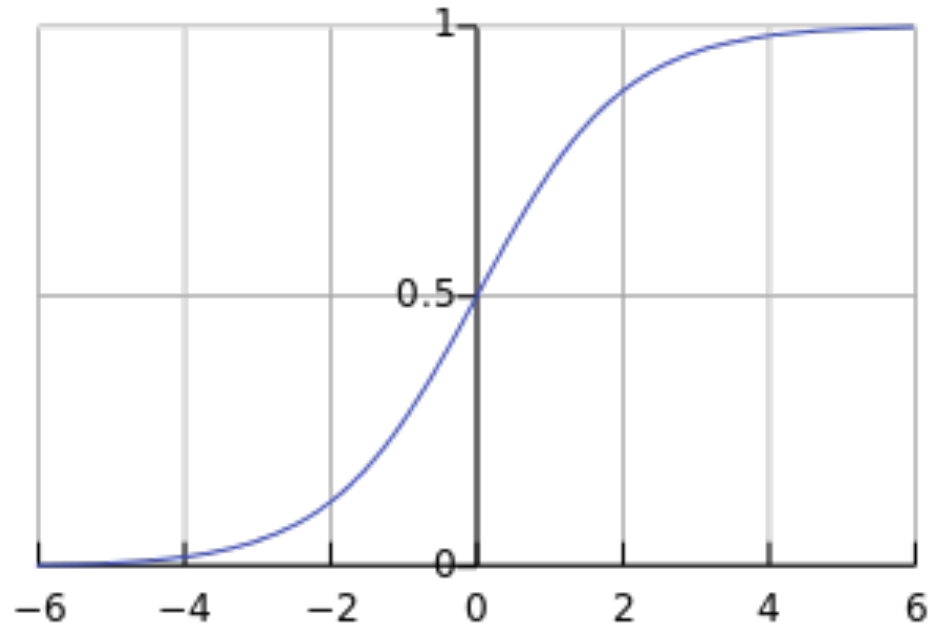$$z \in \mathbb{R} = (-\infty, \infty)$$

$$p(y=1/x) \in [0, 1]$$

Inputs & Outputs

$$z \to \infty$$

$$\left( \frac{1}{1+e^{-z}} \right)$$
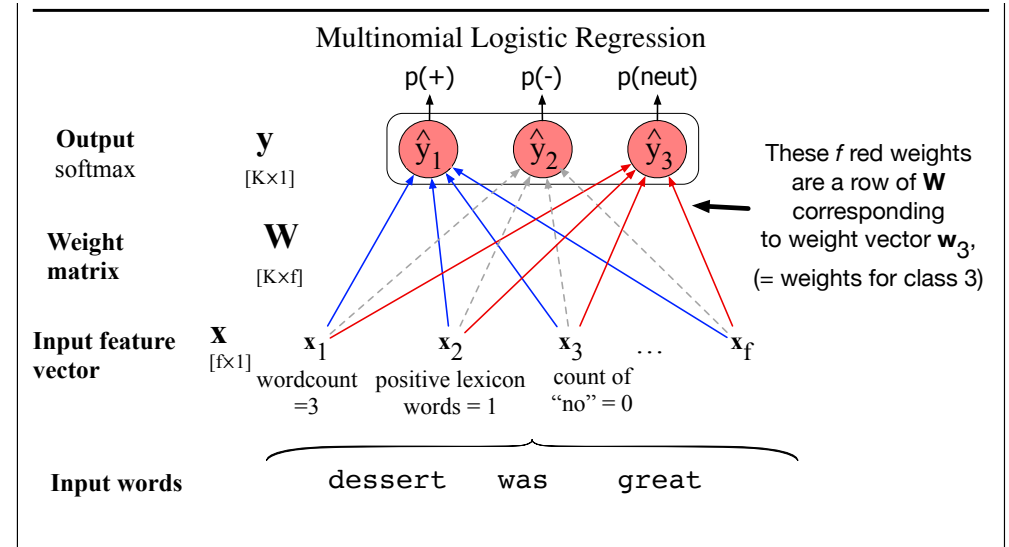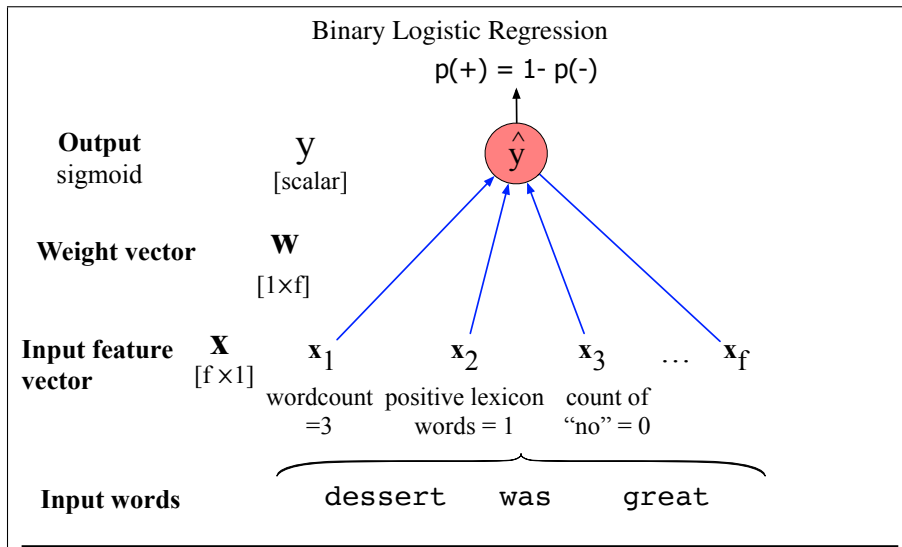
$$\hookrightarrow 1 \xrightarrow{\to 0}$$

# Logistic Function

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

# Multiclass Logistic Regression

*Softmax Reg* is Mult. LR

- Generalize to K>2 classes
- Each class has its own weight vector (across all features: e.g. BOW counts)

**Binary Logistic Regression**

$p(+) = 1 - p(-)$

| | | |
|---|---|---|
| **Output** sigmoid | $y$ [scalar] | |
| **Weight vector** | $\mathbf{W}$ [1×f] | |
| **Input feature vector** | $\mathbf{X}$ [f ×1] | $x_1$ wordcount =3    $x_2$ positive lexicon words = 1    $x_3$ count of "no" = 0    …    $x_f$ |
| **Input words** | | dessert    was    great |

**Multinomial Logistic Regression**

$p(+)$    $p(-)$    $p(neut)$

| | | |
|---|---|---|
| **Output** softmax | $y$ [K×1] | $\hat{y}_1$    $\hat{y}_2$    $\hat{y}_3$ |
| **Weight matrix** | $\mathbf{W}$ [K×f] | |
| **Input feature vector** | $\mathbf{X}$ [f×1] | $x_1$ wordcount =3    $x_2$ positive lexicon words = 1    $x_3$ count of "no" = 0    …    $x_f$ |
| **Input words** | | dessert    was    great |

These *f* red weights are a row of **W** corresponding to weight vector $\mathbf{w}_3$, (= weights for class 3)

---

$p(+)$    $p(-)$    $p(neut)$

| | | |
|---|---|---|
| **Output** softmax | $y$ [K×1] | $\hat{y}_1$    $\hat{y}_2$    $\hat{y}_3$ |
| **Weight matrix** | $\mathbf{W}$ [K×f] | |
| **Input feature vector** | $\mathbf{X}$ [f×1] | $x_1$ wordcount    $x_2$ positive lexicon    $x_3$ count of "no" = 0    …    $x_f$ |

These *f* red weights are a row of **W** corresponding to weight vector $\mathbf{w}_3$, (= weights for class 3)

16

# Multiclass Logistic Regression

- Weight vector for each class

$$\beta_k \in \mathbb{R}^f$$

$$\text{for } k = 1 \dots K$$

$$\text{e.g., } \beta_{k,f} = 3.2$$

- Prediction: dot product for each class

$$z_1 = \beta_1^T x = \sum_{j=1}^{f} \beta_{1,j} x_j$$

$$z_2 = \beta_2^T x$$

$$\vdots$$

$$z_K = \beta_K^T x$$

- Predicted probabilities: apply the *softmax function* to normalize

$$p(y=k \mid x) = \frac{e^{z_k}}{\sum_{\ell=1}^{K} e^{z_\ell}}$$

# Why the softmax function?

want $\left[ p(y=1|x), \; p(y=2|x), \; \ldots \; p(y=k|x) \right]$

① Non-neg       ② Sum to 1

have $\left[ z_1, \; z_2, \; \ldots \; z_k \right]$

$$ e^{z_1} \quad e^{z_2} \quad \ldots \quad e^{z_k} $$

$$ \Rightarrow \quad \frac{e^{z_2}}{\sum_{l=1}^{k} e^{z_l}} $$

# NB as Log-Linear Model

$$P(\text{spam}|D) \propto P(\text{spam}) \cdot \prod_{w_i \in D} P(w_i|\text{spam})$$

# NB as Log-Linear Model

$$P(\text{spam}|D) \propto P(\text{spam}) \cdot \prod_{w_i \in D} P(w_i|\text{spam})$$

$$P(\text{spam}|D) \propto P(\text{spam}) + \prod_{w_i \in \text{Vocab}} \cdot P(w_i|\text{spam})^{x_i}$$

# NB as Log-Linear Model

$$P(\text{spam}|D) \propto P(\text{spam}) \cdot \prod_{w_i \in D} P(w_i|\text{spam})$$

$$P(\text{spam}|D) \propto P(\text{spam}) \cdot \prod_{w_i \in \text{Vocab}} \cdot P(w_i|\text{spam})^{x_i}$$

$$\log[P(\text{spam}|D)] \propto \log[P(\text{spam})] + \sum_{w_i \in \text{Vocab}} x_i \cdot \log[P(w_i|\text{spam})]$$

# NB as log-linear model

$$P(\text{spam} \mid D) = \frac{1}{Z} P(\text{spam}) \prod_{t=1}^{\text{len}(D)} P(w_t \mid \text{spam})$$

$$P(\text{spam} \mid D) = \frac{1}{Z} P(\text{spam}) \prod_{w \in \mathcal{V}} P(w \mid \text{spam})^{x_w}$$

$$\log P(\text{spam} \mid D) = \log P(\text{spam}) + \sum_{w \in \mathcal{V}} x_w \log P(w \mid \text{spam}) - \log Z$$

# NB as Log-Linear Model

In both NB and LogReg

we **compute the dot product!**

# NB vs. LogReg

- Both compute the dot product


- **NB**: sum of log probs; **LogReg**: logistic fun.

# Learning Weights

- **NB**: learn conditional probabilities separately via **counting**

- **LogReg**: learn weights **jointly**

# Learning Weights

- given: a set of **feature vectors** and **labels**

- goal: learn the weights.

# Learning Weights

$$x_{00} \quad x_{01} \quad \ldots \quad x_{0m} \quad y_0$$

$$x_{10} \quad x_{11} \quad \ldots \quad x_{1m} \quad y_1$$

$$\vdots \quad\quad \vdots \quad\quad \ddots \quad\quad \vdots \quad\quad \vdots$$

$$x_{n0} \quad x_{n1} \quad \ldots \quad x_{nm} \quad y_n$$

n examples; xs - features; ys - class

# Learning Weights

We know:

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

$$g(z) = \frac{1}{1 + e^{-z}} \qquad P(y = 1 \mid x) = g\left(\sum_{j=1}^{N_{feat}} \beta_j x_{ij}\right)$$

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

So let's try to maximize probability of the entire dataset - **maximum likelihood estimation**

# Learning Weights

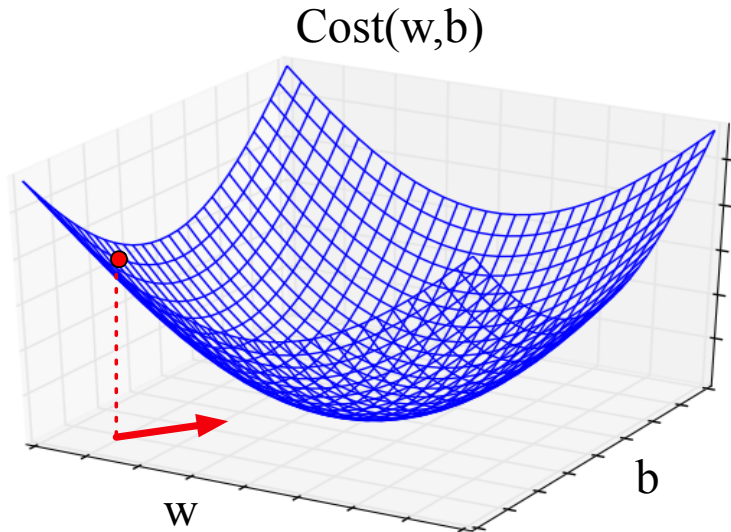So let's try to maximize probability of the entire dataset - **maximum likelihood estimation**

$$\beta^{MLE} = \arg\max_{\beta} \log P(y_0, \ldots, y_n | \mathbf{x_0}, \ldots, \mathbf{x_n}; \beta)$$

# Gradient ascent/descent learning

$$\beta^{MLE} = \arg\max_{\beta} \log P(y_0, \ldots, y_n | \mathbf{x_0}, \ldots, \mathbf{x_n}; \beta)$$

- Follow direction of *steepest ascent*. Iterate: $\quad \beta^{(new)} = \beta^{(old)} + \eta \dfrac{\partial \ell}{\partial \beta}$

Cost(w,b)

$\left( \dfrac{\partial \ell}{\partial \beta_1}, \ldots, \dfrac{\partial \ell}{\partial \beta_J} \right)$ : Gradient vector (vector of per-element derivatives)

GD is a generic method for optimizing differentiable functions — widely used in machine learning!

w

b

# Pros & Cons

- LogReg doesn't assume independence
  - better calibrated probabilities


- NB is faster to train; less likely to overfit

# NB & Log Reg

- Both are linear models:

$$z = \sum_{j=1}^{\text{Nfeat}} \beta_j x_{ij}$$

- Training is different:
  - NB: weights trained independently
  - LogReg: weights trained jointly

# Overfitting and generalization

- Overfitting: your model performs overly optimistically on training set, but generalizes poorly to other data (even from same distribution)
- To diagnose: separate training set vs. test set.
- How did we regularize Naive Bayes and language modeling?


- For logistic regression: L2 regularization for training

# Regularization tradeoffs

- No regularization    <-------------->   Very strong regularization

# Visualizing a classifier in feature space

*"Bias term"*
↓

Feature vector $\quad x = (1, \ \text{count "happy"}, \ \text{count "hello"}, ...)$

Weights/parameters $\quad \beta =$
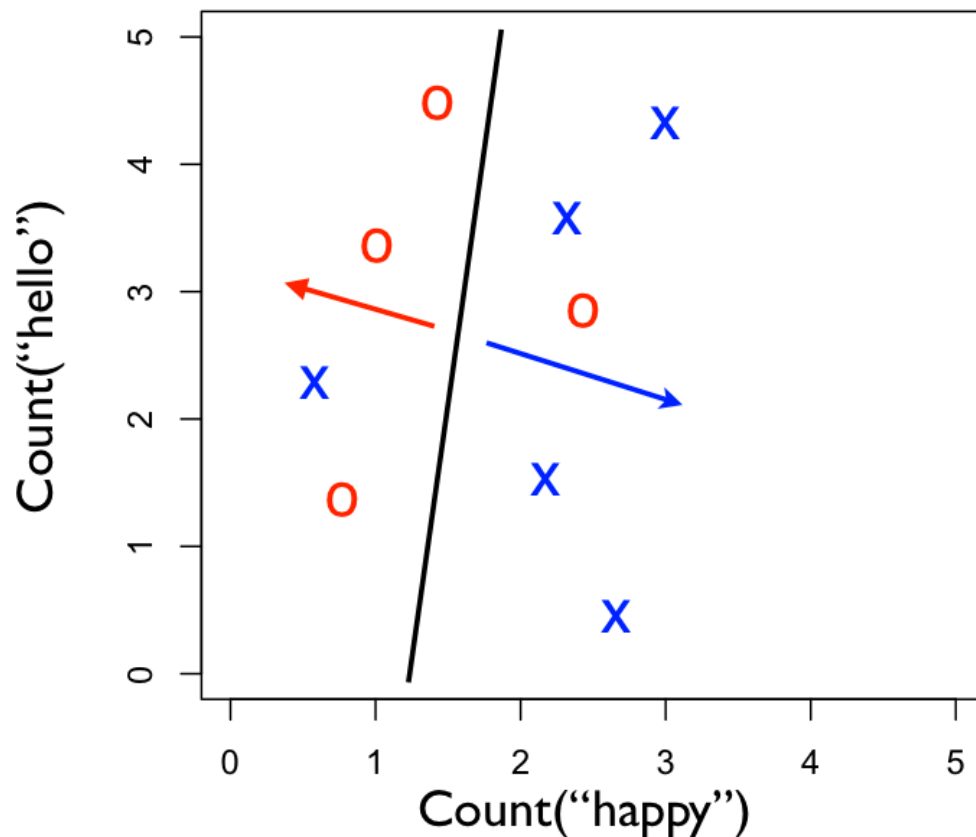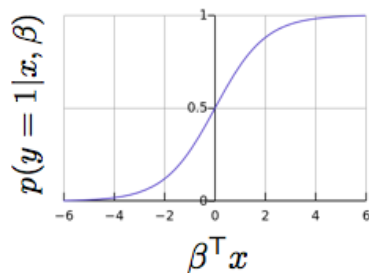
**50% prob where**
$$\beta^{\mathsf{T}} x = 0$$

Predict y=1 when
$$\beta^{\mathsf{T}} x > 0$$

Predict y=0 when
$$\beta^{\mathsf{T}} x \leq 0$$

# Logistic regression wrap-up

- Given you can extract features from your text, logistic regression is the best, easy-to-use, method
  - Logistic regression with BOW features is an excellent baseline method to try at first
  - Will be a foundation for more sophisticated models, later in course
- Always regularize your LR model
- We recommend using the implementation in scikit-learn
  - Useful: CountVectorizer to help make BOW count vectors

- Next: but where do the LABELS in supervised learning come from?