

# Linear Regression, Logistic Regression & Naïve Bayes

04 – ML MODELS



# Agenda Today

- > Linear Regression
  - > Gradient Descent
- > Logistic Regression
- > Naïve Bayes

## Recap: Bias – Variance tradeoff

- > What is the difference between Bias and Variance?
- > What is the difference between under- and over-fitting?
- > How can we overcome these issues?

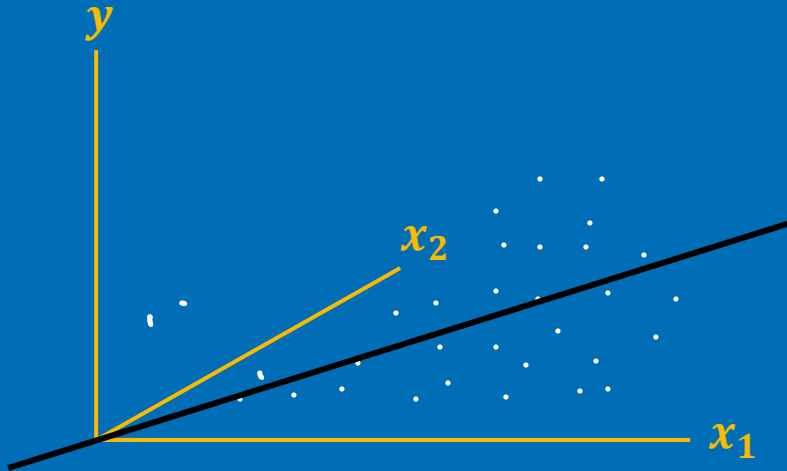
## Recap: Canonical Types // output size // un/supervised?

How would you tackle the following tasks? What *canonical problem type* are they? What is their output space *size*?

- > Assign one of 36 POS tags to each word of a given sentence.
- > Decide if an email is a spam.
- > Determine if a comment is hateful, toxic, threat, insult, or a combination.
- > Automatically extract and assign a set of 10 topics from 1000 documents.
- > Given a comment about a product, predict the score the user gave.

# Linear Regression

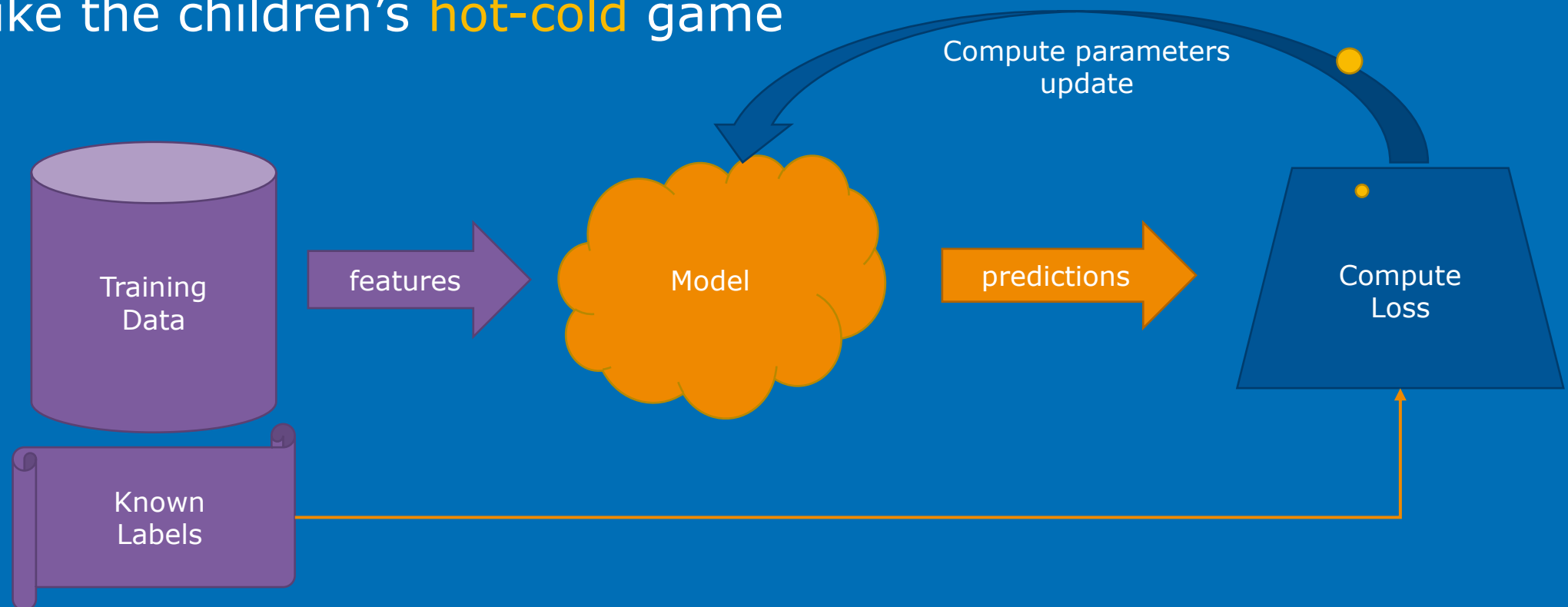
- > A **Supervised Learning** method
- > Predicting a value  $\hat{y} \in \mathbb{R}$  from  $m$  features  $x \in \mathbb{R}^m$  by learning a linear function:  $f(x) = y$
- >  $f$  is of the form:  $\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 \dots = \beta_0 + x^T \beta$



What should/can be  
the loss function?

# Linear Regression: Parameter Estimation

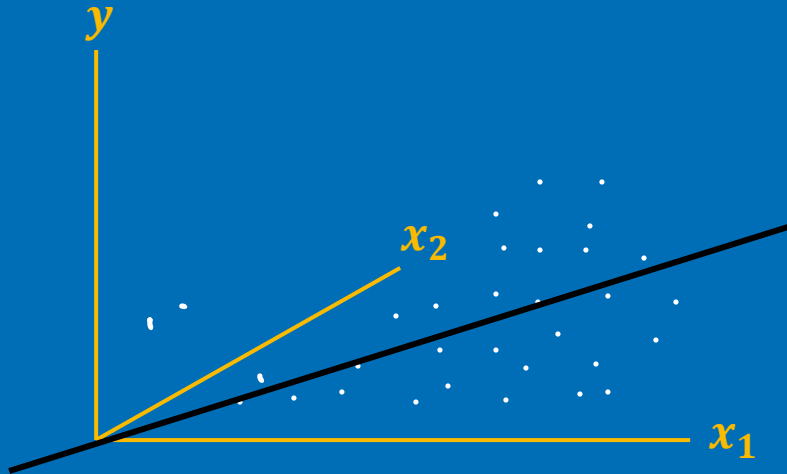
Finding the smallest lost:  
Like the children's **hot-cold** game



# Linear Regression: Parameter Estimation

We can estimate the best parameters  $\langle \beta_0, \beta_1 \rangle$  by minimizing a *loss* function.

A possible *loss*: the **average of squared errors (MSE)** between the **true values** and the **predictions** with these parameters:



$$\hat{\theta} = \arg \min_{\theta = \langle \beta_0, \beta_1 \rangle} \frac{1}{2n} \sum_{i=1}^n \left( y_i - (\beta_0 + x_i^T \beta) \right)^2$$

# Linear Regression: Parameter Estimation

Simplified case with 2 parameters:

$$\hat{y} = \beta_0 + \beta_1 x_1$$

How to find the best parameters  $\langle \beta_0, \beta_1 \rangle$ ?

1. Initialize with random values
2. Calculate the Loss: e.g.,  $\sum_{i=1}^n (y_{true} - y_{pred})^2$
3. Calculate the necessary parameters change.
4. Repeat, until the average loss is very small, or the parameter change is minimal.



# Linear Regression: Parameter Estimation

Simplified case with 2 parameters:

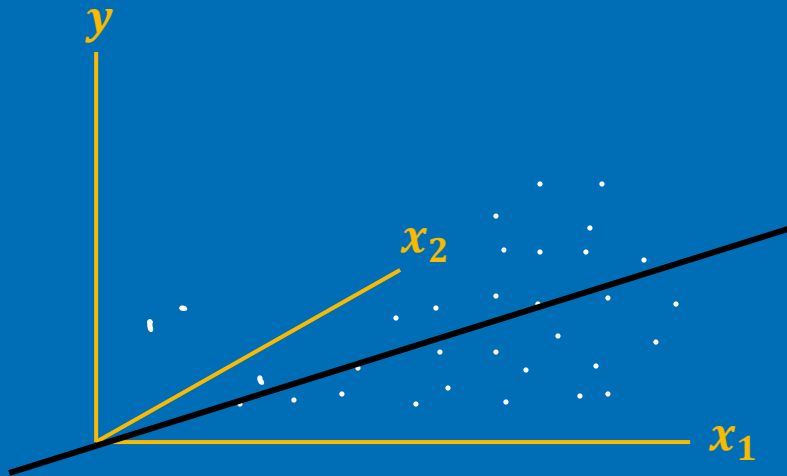
$$\hat{y} = \beta_0 + \beta_1 x_1$$

How to find the best parameters  $\langle \beta_0, \beta_1 \rangle$ ?

1. Initialize with random values
2. Calculate the Loss: e.g.,  $\sum_{i=1}^n (y_{true} - y_{pred})^2$
3. Calculate the necessary parameters change.
4. Repeat, until the average loss is very small, or the parameter change is minimal.

# Linear Regression: Least Mean Squares (LMS)

Parameter Estimation is done through calculating the cost(=loss) function *partial derivatives* with respect to each  $\beta$  and finding its global minima with a method called:  
**Gradient Descent**



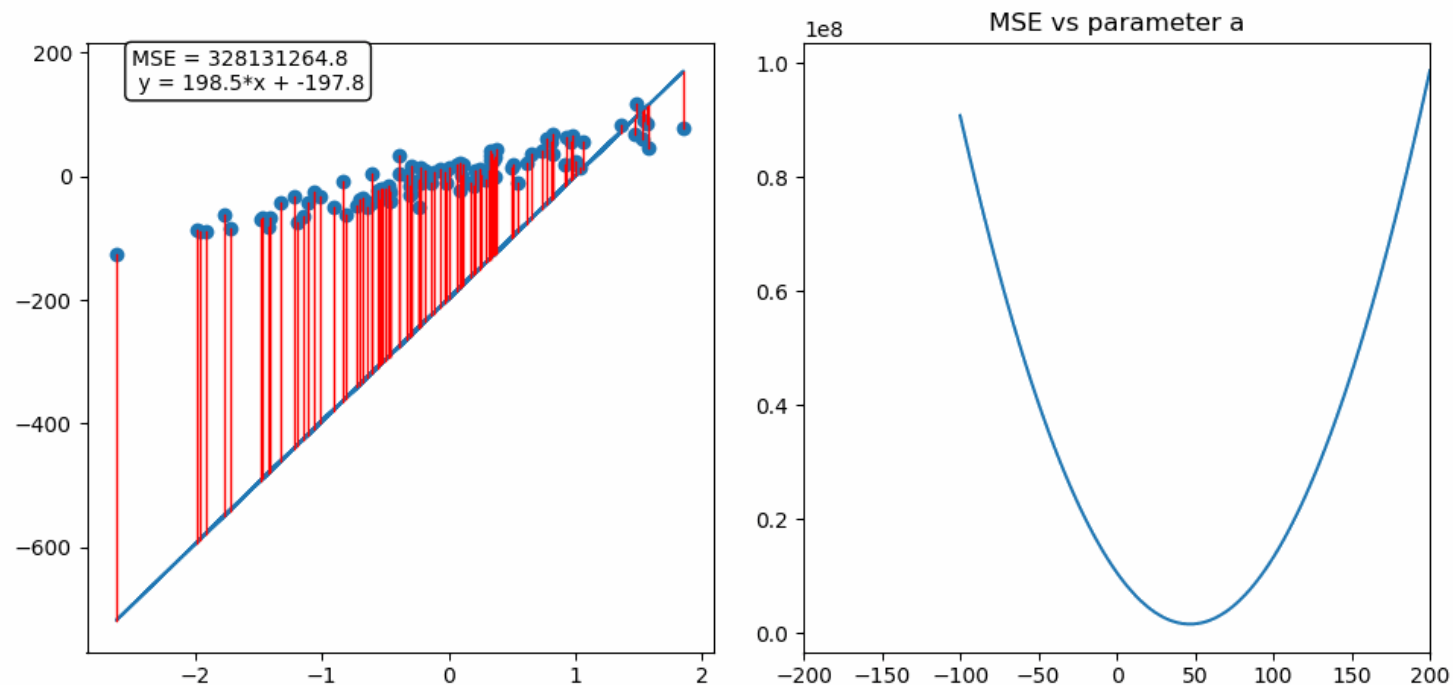
$$\hat{\theta} = \arg \min_{\theta = \langle \beta_0, \beta_1 \rangle} \frac{1}{2n} \sum_{i=1}^n \left( y_i - (\beta_0 + x_i^T \beta) \right)^2$$

# Optimizing Linear Regression with Gradient Descent

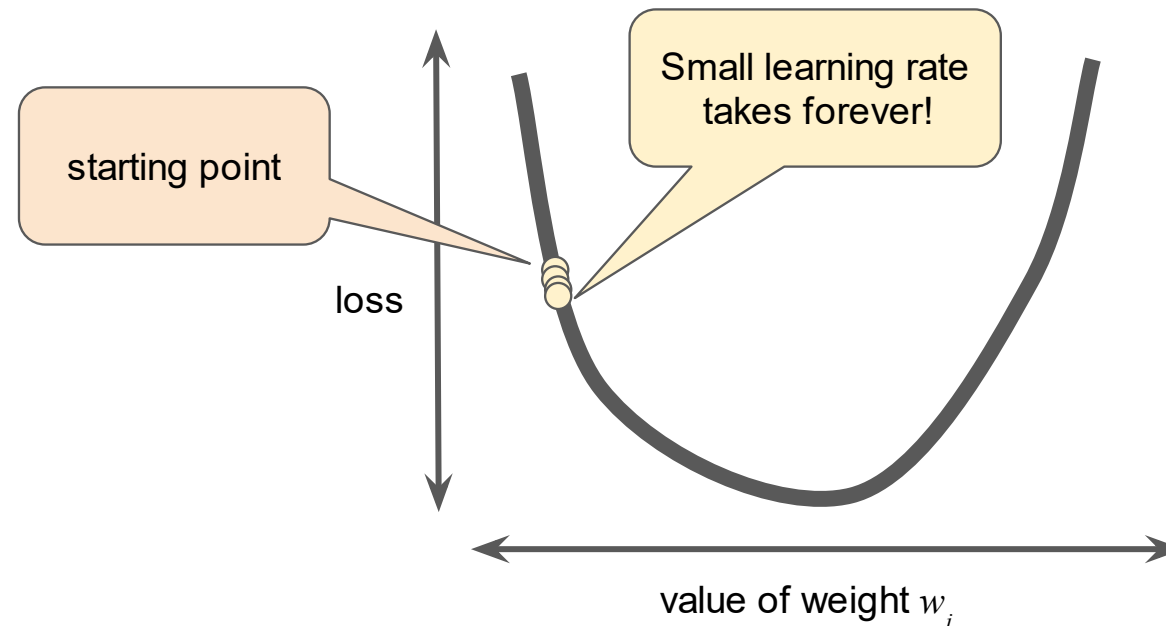
- >  $\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 \dots$
- > Calculating partial derivatives (**gradient**)
- > Finding the slope – which direction minimizes (**descent**) the gradients?  
And by how much?  
(Multiplied by some fraction - **Learning Rate**)

# Gradient Descent

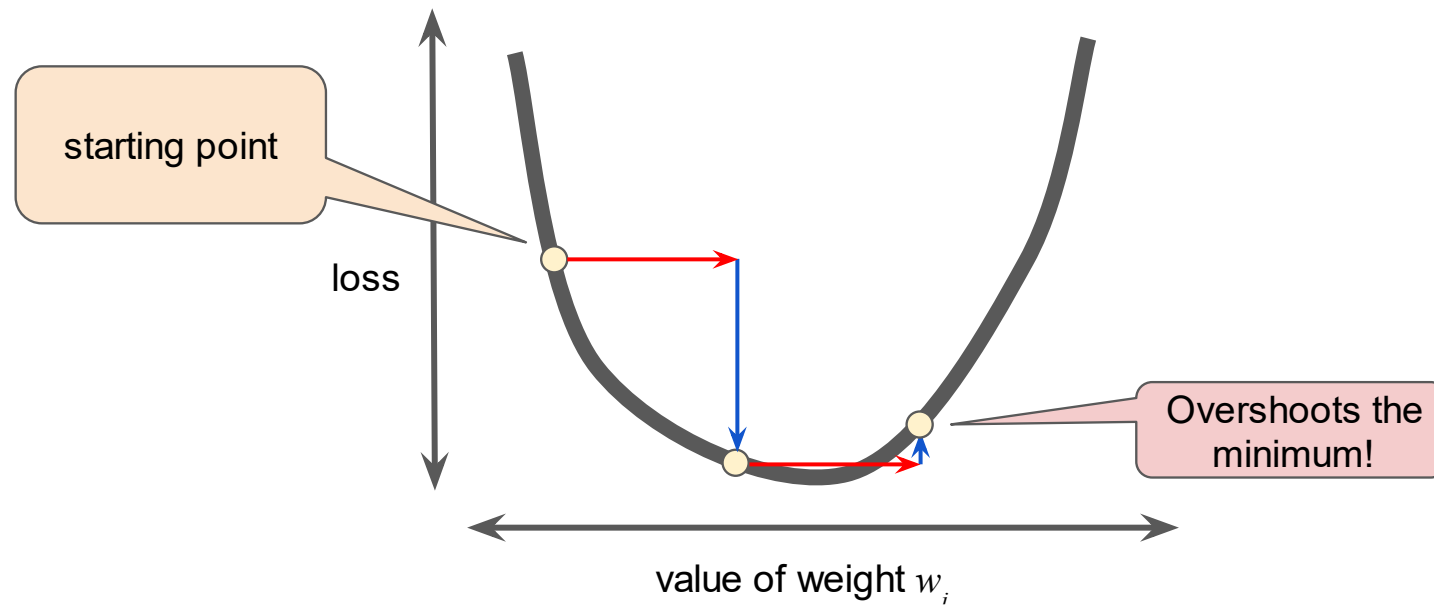
- > Minimizing the partial derivatives (like in high school) – to find the global minimum
- > Often is multiplied by some **learning-rate**
  - if too big – jumps too far
  - if too small – takes too long



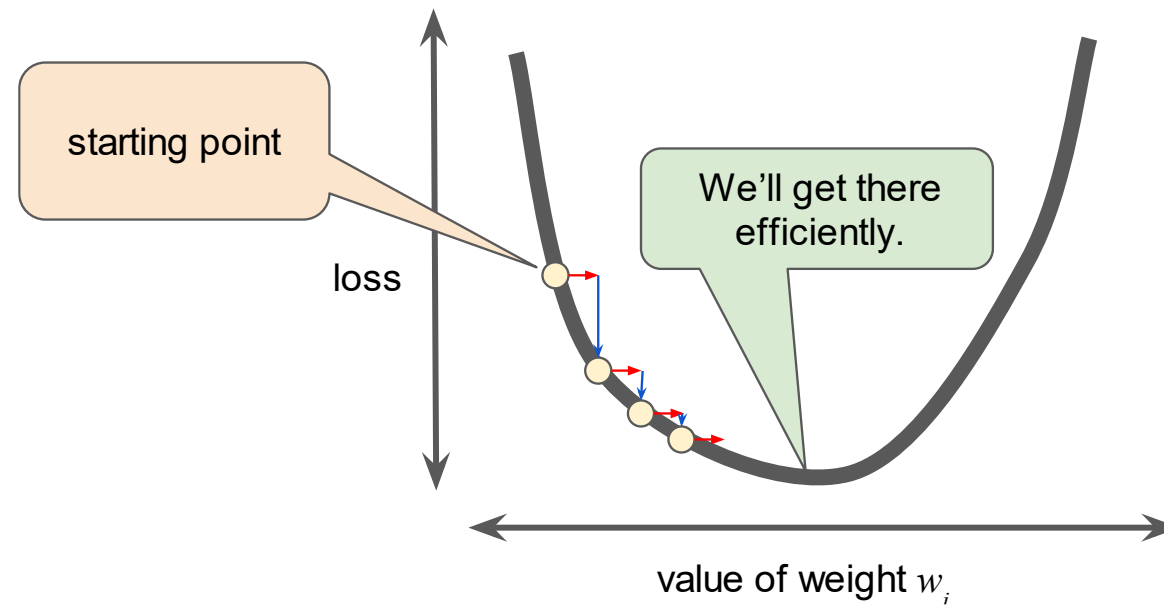
# Learning Rate – too small



# Learning Rate – too large

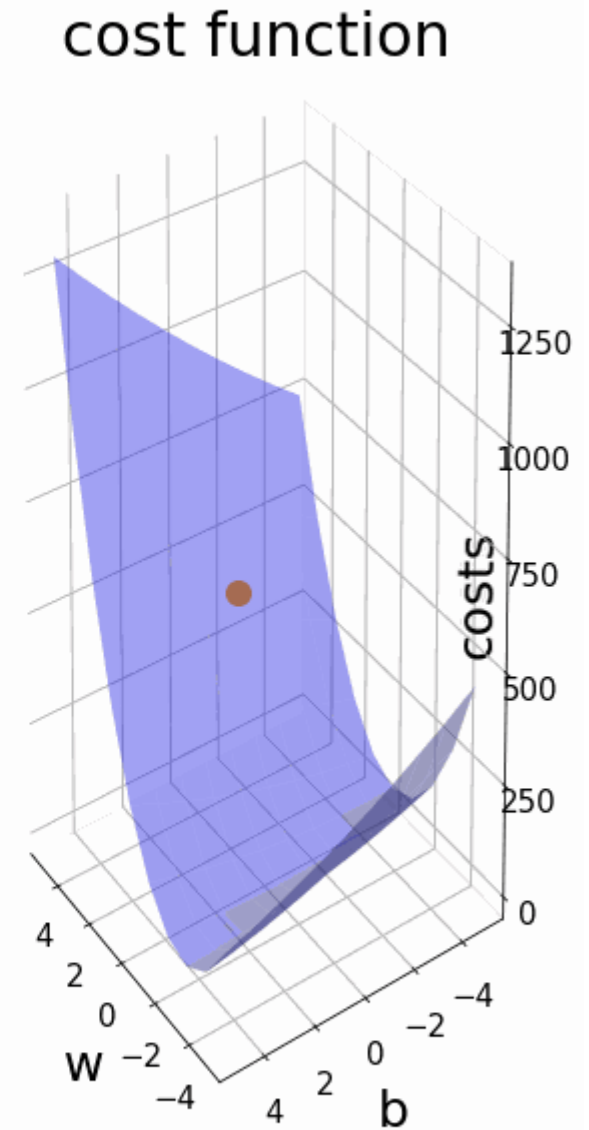
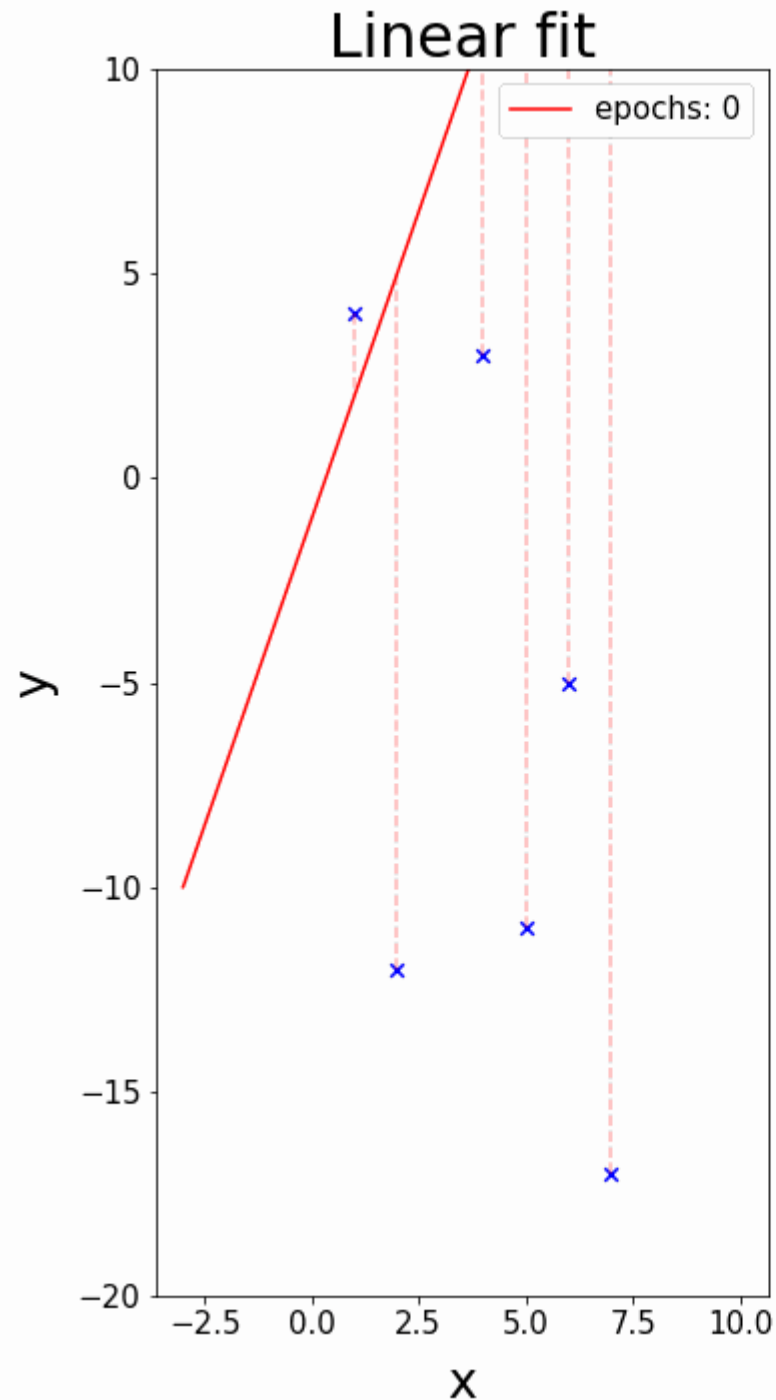


# Learning Rate – just right



## Two-dimensional

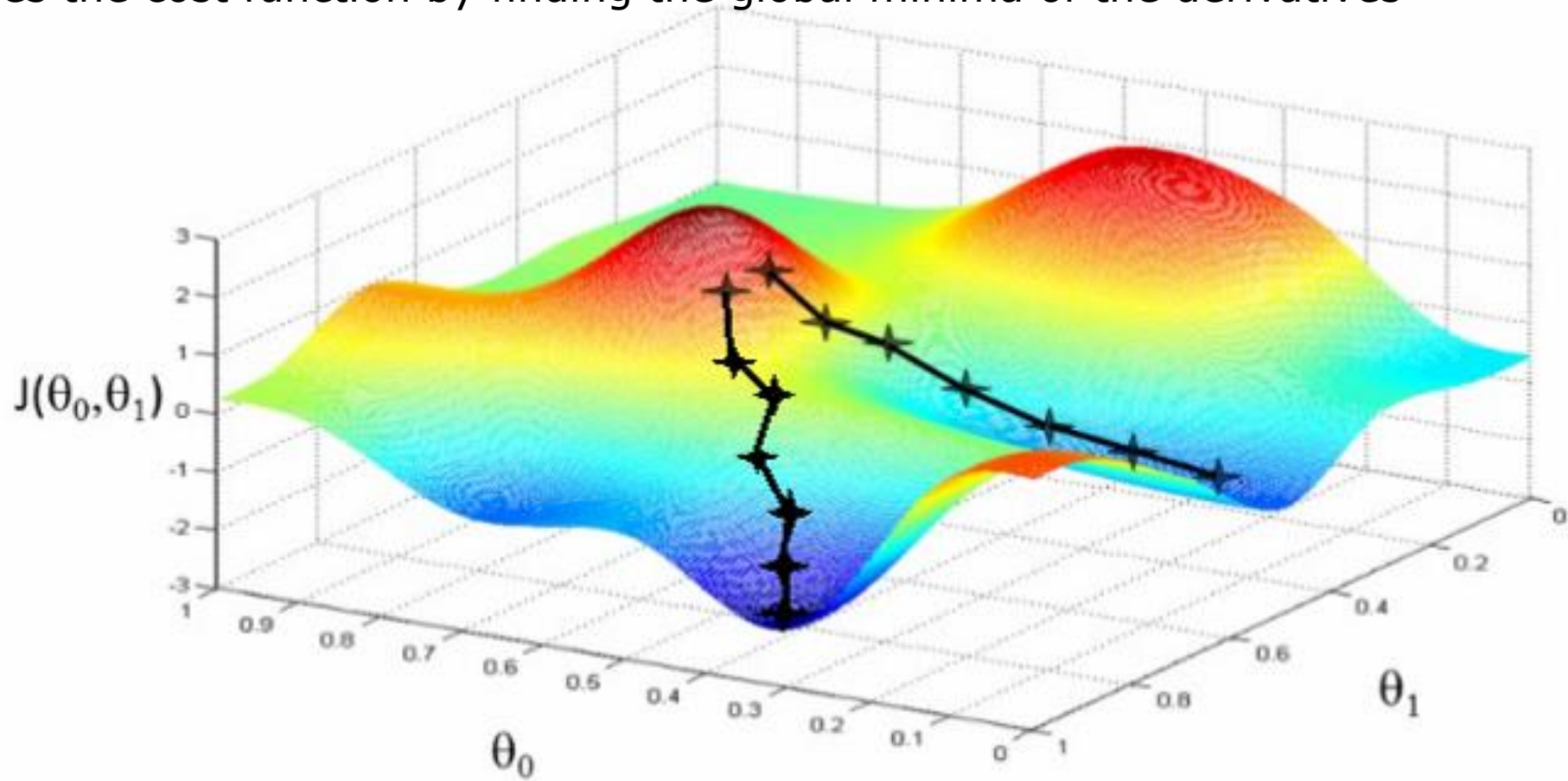
- >  $\hat{y} = \beta_0 + \beta_1 x$
- > Minimizing both -  
 $\beta_0$  &  $\beta_1$



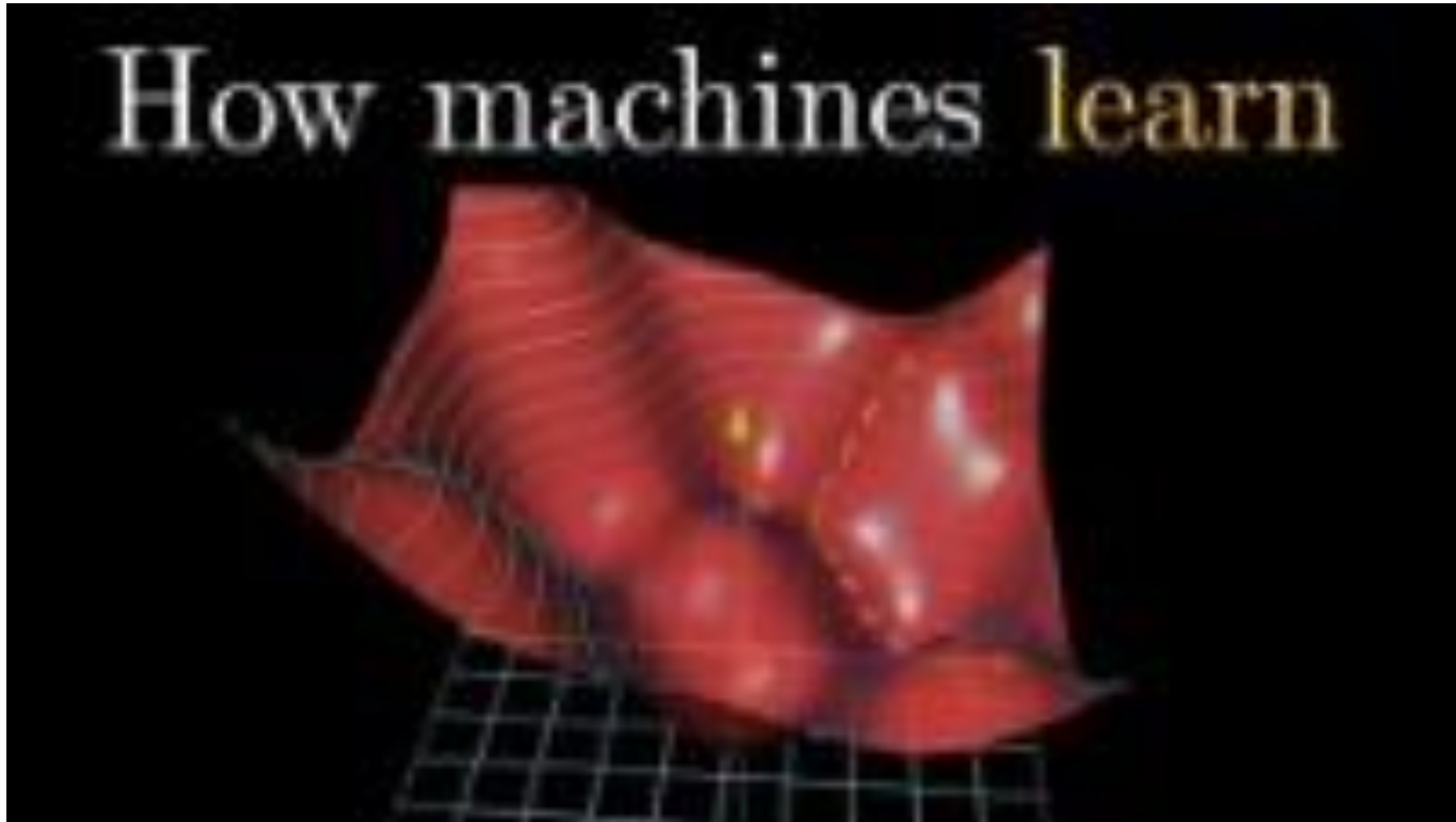


# Gradient Descent

Minimizes the cost function by finding the global minima of the derivatives



# Gradient Descent: Recommended Video



## Back to Linear Regression

- > Calculates the best linear fit of the given features
- > Assumes the features are **independent (!)**
- > Requires numerical features
  - > Which features would you craft for the following problem:  
Given an essay, predict the student's score (0-100)

# Linear Regression - Regularization

- > ML training can cause overfitting.
  - > Q: What would overfitting look like for **Linear** regression?
- > How can we prevent it?  
Regularization!

$$\hat{\theta} = \arg \min_{\theta = \langle \beta_0, \beta_1 \rangle} \frac{1}{2n} \sum_{i=1}^n \left( y_i - (\beta_0 + x_i^T \beta) \right)^2$$

# Linear Regression - Regularization

- > ML training can cause overfitting.
  - > What would overfitting look like for **Linear** regression?
- > How can we prevent it?

$L_2$  Regularization: Euclidean Distance

$$\hat{\theta} = \arg \min_{\theta = \langle \beta_0, \beta_1 \rangle} \frac{1}{2n} \sum_{i=1}^n \left( y_i - (\beta_0 + x_i^T \beta) \right)^2 + \frac{\lambda}{2} \sqrt{\sum_{j=1}^m \beta_j^2}$$

# Linear Regression - Regularization

- > ML training can cause overfitting.
  - > What would overfitting look like for **Linear** regression?
- > How can we prevent it?

$L_1$  Regularization: Manhattan (Taxicab) Distance

$$\hat{\theta} = \arg \min_{\theta = \langle \beta_0, \beta_1 \rangle} \frac{1}{2n} \sum_{i=1}^n \left( y_i - (\beta_0 + x_i^T \beta) \right)^2 + \lambda \sum_{j=1}^m |\beta_j|$$

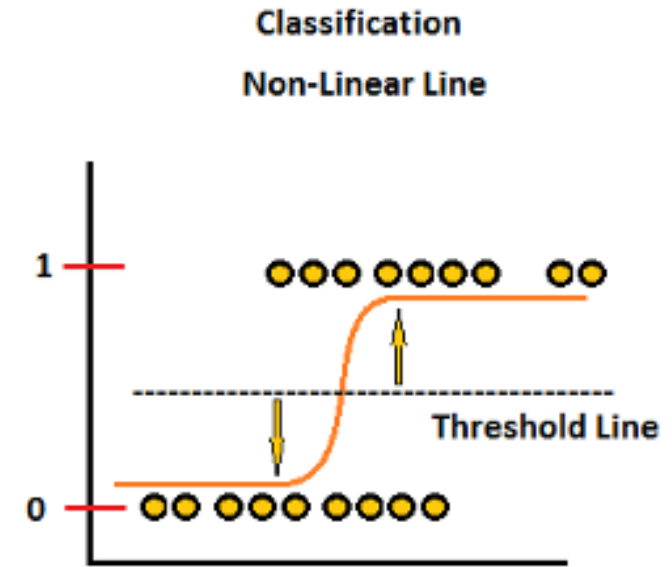
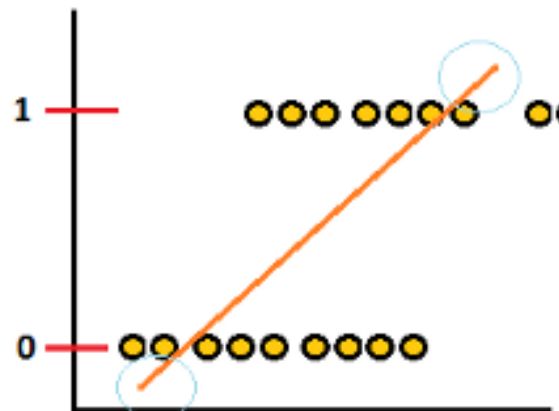
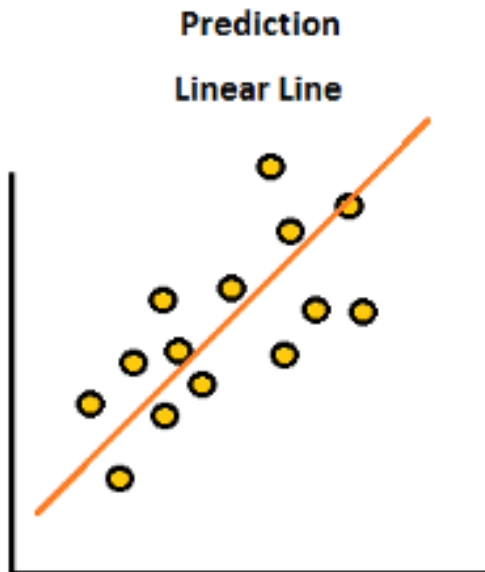
# **From Linear to Logistic**

Classification using  
Logistic Regression



# Binary Logistic Regression

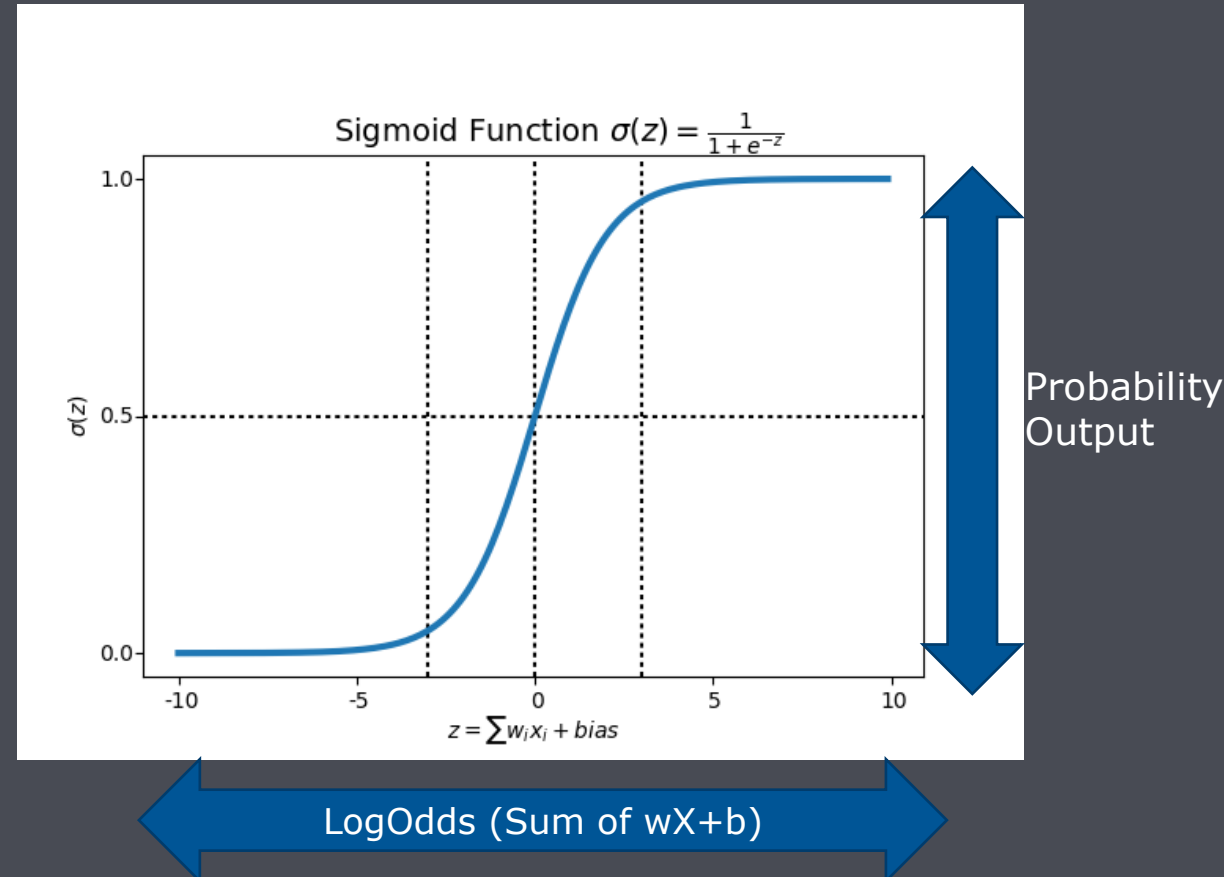
- > Linear Regression: Predict an essay's **score** (e.g., readability)
- > Logistic Regression: Predict a **probability** (Readable or not)





# Linear Logistic Regression

- > Gives a **probability** estimation
- > By wrapping the **linear function** in a non-linear function.  
E.g., **sigmoid function**:
$$y' = \frac{1}{1 + e^{-(w^T x + b)}}$$
- > Squishing the linear output through sigmoid



# Linear Logistic Regression – Sigmoid Function

- > Probability vs Odds.
- > Probability: 0...1
- > Odds:  $\frac{P}{1-P}$
- > Logistic Function:

$Y$	$1$	$0$
$Pr(Y=1)$	$P$	$1-P$

*\*P= Success, 1-P = Failure*

$$\ln\left(\frac{P}{1-P}\right) = \beta_0 + \beta_1 x$$

$$\frac{P}{1-P} = e^{\beta_0 + \beta_1 x}$$

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

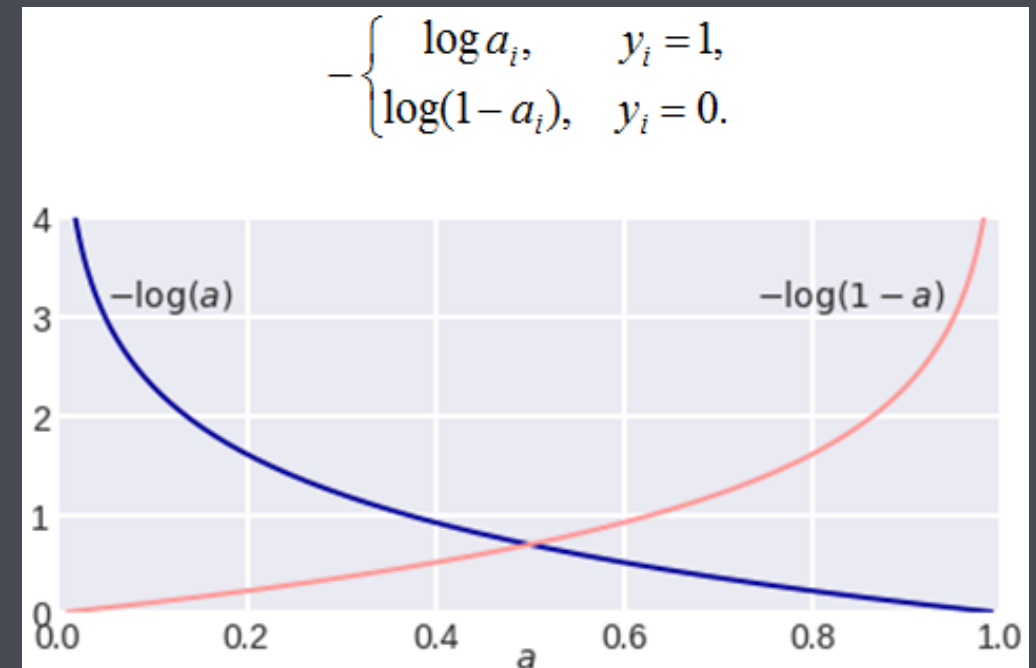
# Linear Logistic Regression - Loss Function

- > In Linear Regression we used **Mean Squared Error (MSE)**.
- > In Logistic Regression we need a different loss:

$$\text{Log Loss} = \sum_{(x,y) \in D} -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

- > As we reach closer to the side of the bars, the probability gets high quickly.

$(x, y) \in D$  – data and its labels  
 $y$  – the true label (between 0 and 1)  
 $y'$  – the predicted label (between 0 and 1)



# Regularization

- > Regularization is very important in logistic regression modeling.
- > Without it, the loss would go towards 0 in high dimensions (due to the asymptotic nature of logistic regression)
- > Strategies to dampen model complexity:
  1.  $L_2$  regularization.
  2. Early stopping

## Example: Sentiment Classification of movie reviews

- > Problem definition: Estimate if a review is positive or not.
- > Features:
  1. `count(positive lexicon words ∈ doc)`
  2. `count(negative lexicon words ∈ doc)`
  3. `{ 1 – “no” found; 0 – Otherwise }`
  4. `count(1st and 2nd pronouns ∈ doc)`
  5. `{ 1 – “!” found; 0 – Otherwise }`
  6. `log(word count of doc)`

## Example: Sentiment Classification of movie reviews

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 .

## Example: Sentiment Classification of movie reviews

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** .

- |  |             |
|--|-------------|
| 1. count(positive lexicon words ∈ doc)                       | 3           |
| 2. count(negative lexicon words ∈ doc)                       | 2           |
| 3. { 1 – “no” found; 0 – Otherwise }                         | 1           |
| 4. count(1 <sup>st</sup> and 2 <sup>nd</sup> pronouns ∈ doc) | 3           |
| 5. { 1 – “!” found; 0 – Otherwise }                          | 0           |
| 6. log(word count of doc)                                    | ln(66)=4.19 |

# Example: Sentiment Classification of movie reviews

Assuming we have a trained model with the following weights:

$$\beta_0 = 0.1, \quad \beta = \begin{matrix} 2.5 & 1. & 3 \\ -5.0 & 2. & 2 \\ -1.2 & 3. & 1 \\ 0.5 & 4. & 3 \\ 2.0 & 5. & 0 \\ 0.7 & 6. & \ln(66)=4.19 \end{matrix}$$

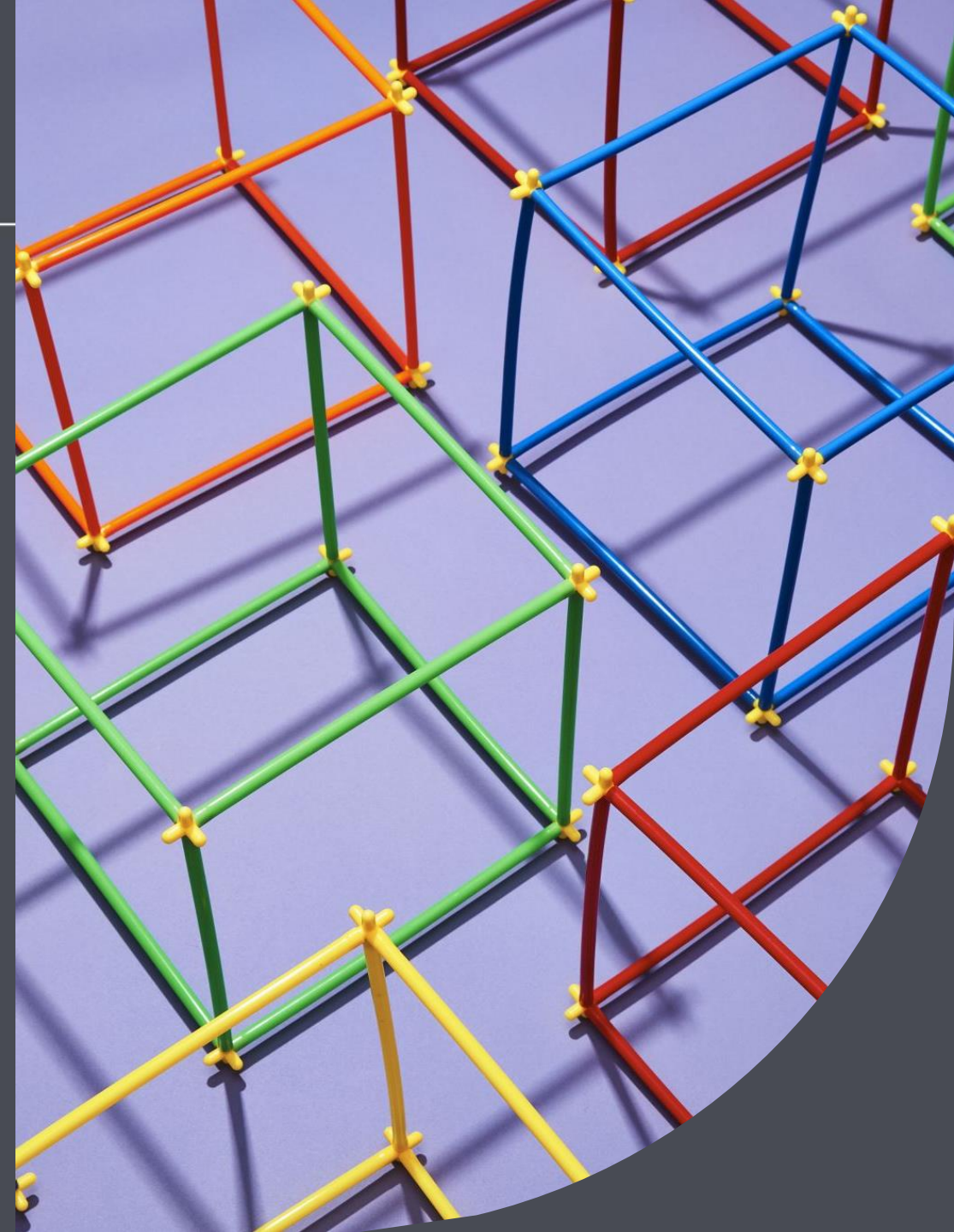
$$\sigma = \frac{1}{1 + e^{-(\beta_0 + \beta x)}}$$

$$P(y = 1|x) = \sigma(\beta \cdot x + \beta_0) = \sigma(.833) = 0.70$$



# Demo

- > [A Neural Network Playground \(tensorflow.org\)](https://www.tensorflow.org/playground)



## What about multiclass?

- > For  $k$  classes, we need  $k$  models
- >  $y$  will be a 1-hot vector:  $[0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0]$
- > Multinomial Logistic Regression

Softmax instead of sigmoid:

$$\text{softmax}(z_i) = \frac{\exp(z_i)}{\sum_{j=1}^k \exp(z_j)} = \frac{\exp(w_k \cdot x + b_k)}{\sum_{j=1}^k \exp(w_j \cdot x + b_j)}$$

# What about multiclass?

Cross entropy loss:

$$-y \cdot \log(p) + (1 - y) \log(1 - p)$$

Which generalizes to:

$$-\sum_{c=1}^K y_c \log(p_c)$$

Softmax instead of sigmoid:

$$\text{softmax}(z_i) = \frac{\exp(z_i)}{\sum_{j=1}^k \exp(z_j)} = \frac{\exp(w_k \cdot x + b_k)}{\sum_{j=1}^k \exp(w_j \cdot x + b_j)}$$

# Understand the math

> Demo

# Linear Logistic Regression – Final Notes

- > Logistic Regression:  
Linear regression wrapped with a non-linear function:
  - > Sigmoid for binary classification
  - > Softmax for multiclass classification
- > Very fast to train
- > Efficient in RAM consumption
- > The basis of Neural Networks

## Recommended Read & Watch for Deeper Dive

- > [Speech and Language Processing. Daniel Jurafsky & James H. Martin. Chapter 5 – Logistic Regression.pdf](#)
- > [Logistic Regression -- Why sigmoid function? \(sebastianraschka.com\)](#)
- > [\(17\) Softmax Function Explained In Depth with 3D Visuals - YouTube](#)
- > [https://ml-cheatsheet.readthedocs.io/en/latest/logistic\\_regression.html](#)
- > [Loss Functions — ML Glossary documentation \(ml-cheatsheet.readthedocs.io\)](#)

# **Naïve Bayes**

Classification with pure Bayesian statistics



# Recap: Maximum Likelihood Estimation (MLE)

> Recall your MLE with Christian



# Naïve Bayes – an Introduction

- > Similar ideas to MLE
- > A probabilistic **classifier**
  - > (a **generative** one)
- > Makes a naïve assumption about the features
- > Assumes feature independence
- > Uses the **Bayes Rule**:

$$P(x|y) = \frac{P(y|x)P(x)}{P(y)}$$

# Bayesian vs Frequentists

## Bayesians:

- Conditional Probability
- Belief-based
- A parameter is a random variable
- Handle uncertainty with opinions
- Goal: Update opinions

## Frequentists:

- Based on data
- Parameters are static
- Probability =  
How many times am I right?
- Uncertainty requires data collection (null hypothesis)
- Goal: decide which action to take

# Bayesian vs Frequentists



# Classification using Naïve Bayes

> Given a document, classify its author:

*p(Mary Shelly | "The fallen angel becomes a malignant devil")*

$$\hat{c} = \arg \max P(c|d) = \arg \max \frac{P(d|c)P(c)}{P(d)}$$

Since our document is the same for all classes, it can be simplified to:

$$\hat{c} = \arg \max P(c|d) = \arg \max P(d|c)P(c)$$

# Classification using Naïve Bayes

> Given a document, classify its author:

*p(Mary Shelly | "The fallen angel becomes a malignant devil")*

$$\hat{c} = \arg \max P(c|d) = \arg \max \frac{P(d|c)P(c)}{P(d)}$$

Since our document is the same for all classes, it can be simplified to:

$$\hat{c} = \arg \max P(c|d) = \arg \max \underbrace{P(d|c)}_{\text{Likelihood}} \underbrace{P(c)}_{\text{Prior}}$$

# Classification using Naïve Bayes

> Given a document, classify its author:

*$p(\text{Mary Shelly} \mid \text{"The fallen angel becomes a malignant devil"})$*

$$\hat{c} = \arg \max P(c|d) = \arg \max P(\text{sentence} \uparrow) P(\text{Mary Shelly})$$

Since our document is the same for all classes, it can be simplified to:

$$\hat{c} = \arg \max P(c|d) = \arg \max \underbrace{P(d|c)}_{\text{Likelihood}} \underbrace{P(c)}_{\text{Prior}}$$

How are these  
probabilities  
obtained?

# Classification using Naïve Bayes

> Given a document, classify its author:

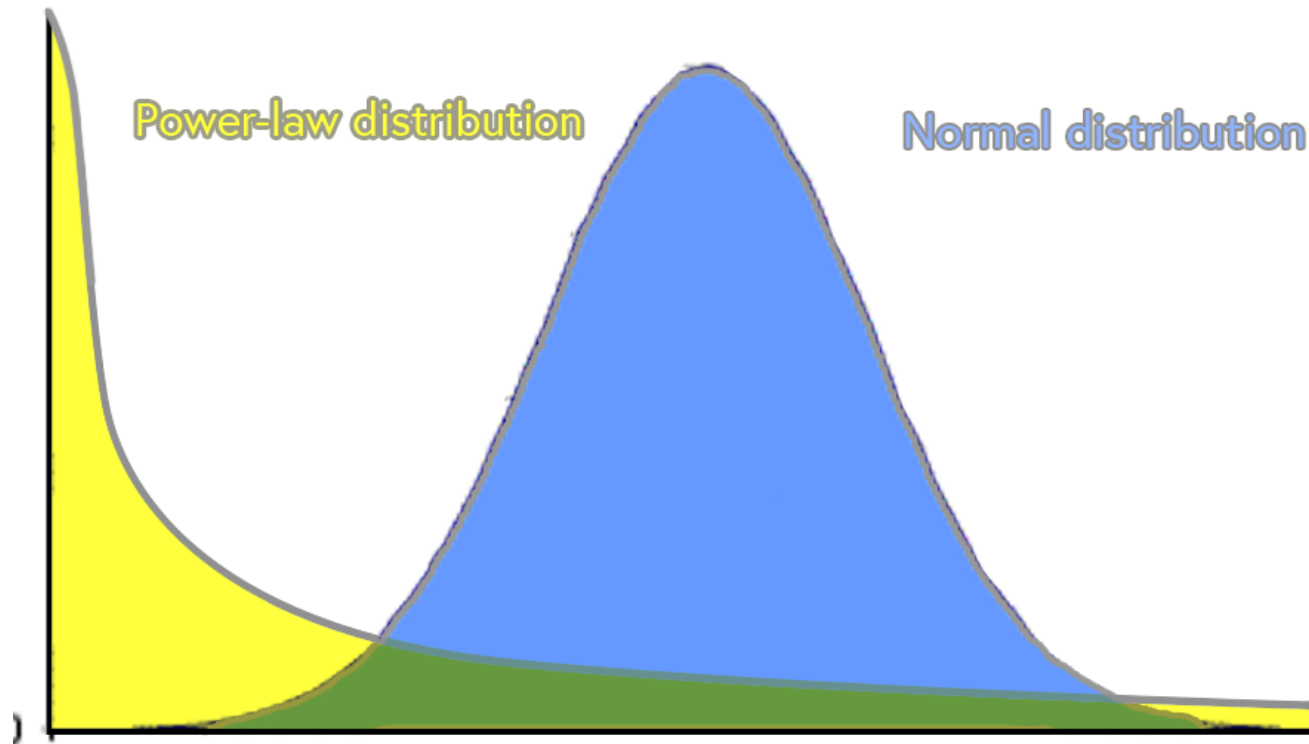
*$p(\text{Mary Shelly} \mid \text{"The fallen angel becomes a malignant devil"})$*

$P(\text{sentence} \uparrow)P(\text{Mary Shelly})$  •  
 $P(\text{sentence} \uparrow)P(\text{Edgar Ellen Poe})$   
 $P(\text{sentence} \uparrow)P(\text{HP Lovecraft})$

Since our document is the same for all classes, it can be simplified to:

$$\hat{c} = \arg \max P(c|d) = \arg \max \underbrace{P(d|c)}_{\text{Likelihood}} \underbrace{P(c)}_{\text{Prior}}$$

# Language data properties



> Zipf Law

> Word frequencies follow a "power-law distribution":

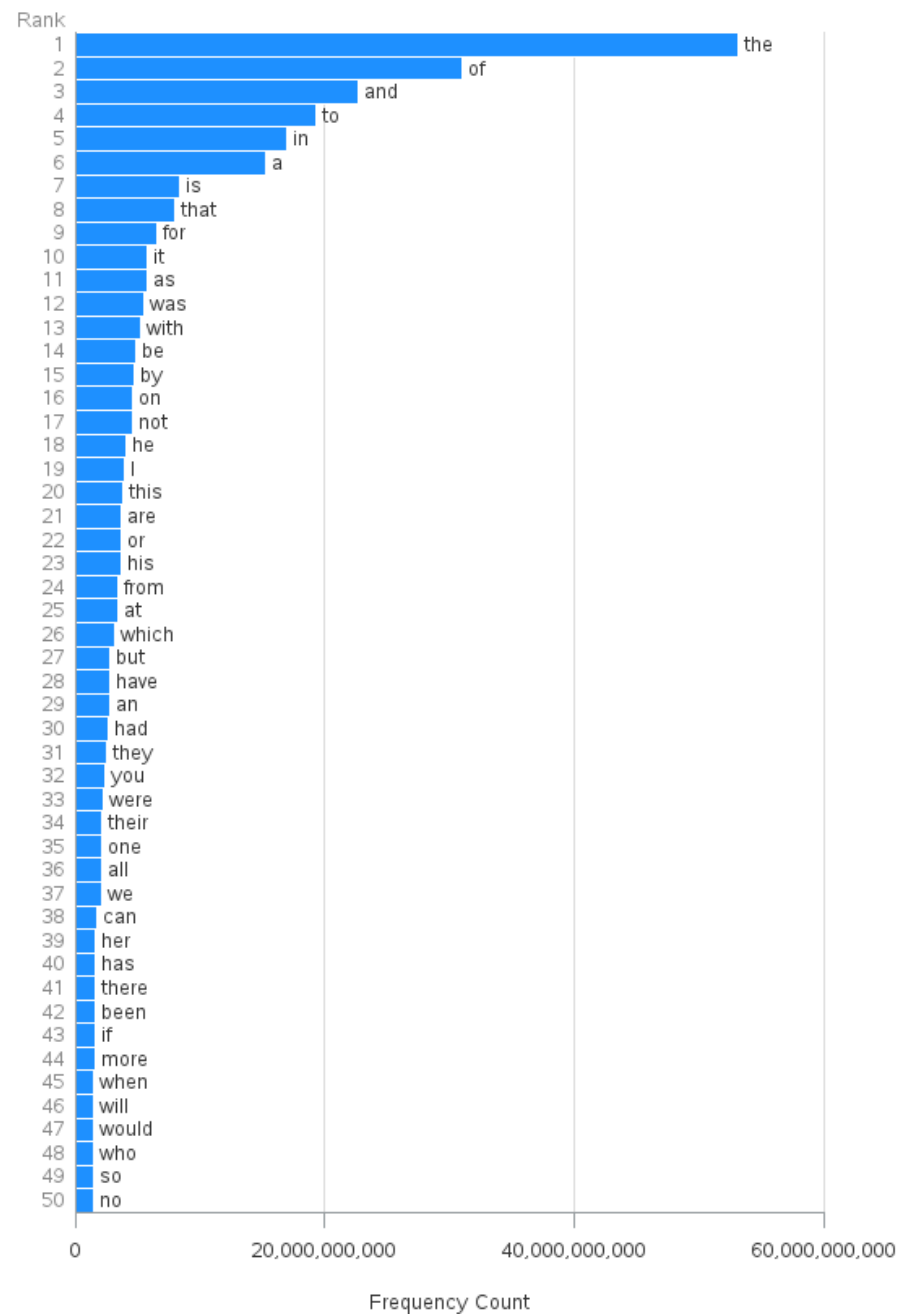
> **Long tail**

> Most events rarely occur



# 50 Most Frequent Words in English Writing

Based on Google books data



## Zipf Law

frequency of a word is inversely proportional to its rank in the frequency table

$$n(r) \propto \frac{1}{r^z} \quad z \approx 1$$

# Zipf Law

In a 43M words text, there are:

- 316,710 unique words (types)
- 144,999 words occur only once
- 42,525 words occur 2 times
- 21,618 words occur 3 times
- 13,306 words occur 4 times
- 9,488 words occur 5 times
- 26,024 words appear >50 times

# Zipf Law

No matter how large the training corpus is:

- > It's likely to contain previously unseen word forms
- > There will be many previously unseen word-*pairs*
- > There will be even more previously unseen word-*triplets*
- > There will be even more previously unseen *sentences*

# Representing text as Features

"to be or not to be"

**Bag of Words** (BoW) - word counts:

{ "be": 2, "to": 2, "not": 1, "or": 1, "something": 0, "else": 0 }

(0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, ... , 1, 0, 1, 0, .... , 0 )

('a', 'the', 'there', 'to', ..., 'be', 'if', 'an' ... , 'or', 'as', 'not'..., 'zebra' )

# Representing text as Features

**"to be or not to be"**

**Bag of Words** (BoW) - word counts:

{ "be": 2, "to": 2, "not": 1, "or": 1, "something": 0, "else": 0 }

**Reweighting:**

TF/IDF or Pointwise Mutual Information – PMI

{ "be": 1.2, "to": 0.1, "not": 0.9, "or": 0.3 }

# TF / IDF + PMI

{ "be": 2, "to": 2, "not": 1, "or": 1 }

$$w_{i,j} = tf_{i,j} \times \log \left( \frac{N}{df_i} \right)$$

$tf_{ij}$  = number of occurrences of  $i$  in  $j$   
 $df_i$  = number of documents containing  $i$   
 $N$  = total number of documents

$$PMI(a, b) = \log \left( \frac{P(a, b)}{P(a)P(b)} \right)$$

# Classification using Naïve Bayes

> Given a document, classify its author:

$p(\text{Mary Shelly} \mid \text{"The fallen angel becomes a malignant devil"})$

$$P(\text{sentence} \uparrow) P(\text{Mary Shelly})$$

(Naively) Assuming feature independence:

$$\begin{aligned} P(\text{sentence}) = & \\ & P(\text{The}) \times P(\text{"fallen"}) \times \\ & P(\text{"angel"}) \times P(\text{"becomes"}) \times P(\text{"a"}) \times \\ & P(\text{"malignant"}) \times P(\text{"devil"}) \end{aligned}$$

# Classification using Naïve Bayes

> Given a document, classify its author:

$p(\text{Mary Shelly} \mid \text{"The fallen angel becomes a malignant devil"})$

$$P(\text{sentence} \uparrow) P(\text{Mary Shelly})$$

(Naively) Assuming feature independence:

$$\begin{aligned} P(\text{sentence} \mid \text{Mary Shelly}) = & \\ & P(\text{"The"} \mid \text{Mary Shelly}) \times P(\text{"fallen"} \mid \text{Mary Shelly}) \times \\ & P(\text{"angel"} \mid \text{Mary Shelly}) \times P(\text{"becomes"} \mid \text{Mary Shelly}) \times P(\text{"a"} \mid \text{Mary Shelly}) \times \\ & P(\text{"malignant"} \mid \text{Mary Shelly}) \times P(\text{"devil"} \mid \text{Mary Shelly}) = \prod_{w \in \text{sentence}} P(w \mid \text{Mary Shelly}) \end{aligned}$$



What happens if a word doesn't appear in the text at all?

# Classification using Naïve Bayes

> Given a document, classify its author:

$p(\text{Mary Shelly} \mid \text{"The fallen angel becomes a malignant devil"})$

$$P(\text{sentence} \uparrow) P(\text{Mary Shelly})$$

(Naively) Assuming feature independence:

Obtaining an **a-priori** by counting word occurrences per class.

$$P(\text{"devil"} \mid \text{Mary Shelly}) = \frac{5}{1000}$$

• "devil" occurrences  
Total words in Mary Shelly documents

# Classification using Naïve Bayes - Unseen Words

> Given a document, classify its author:

$p(\text{Mary Shelly} \mid \text{"The fallen angel becomes a malignant devil"})$

$$P(\text{sentence} \uparrow) P(\text{Mary Shelly})$$

(Naively) Assuming feature independence:

$$P(\text{sentence} \mid \text{Mary Shelly}) =$$

$$P(\text{"The"} \mid \text{Mary Shelly}) \times P(\text{"fallen"} \mid \text{Mary Shelly}) \times$$

$$P(\text{"angel"} \mid \text{Mary Shelly}) \times P(\text{"becomes"} \mid \text{Mary Shelly}) \times P(\text{"a"} \mid \text{Mary Shelly}) \times$$

$$P(\text{"malignant"} \mid \text{Mary Shelly}) \times P(\text{"devil"} \mid \text{Mary Shelly}) \Rightarrow 0$$

$$\frac{0}{1000} = 0$$

# Classification using Naïve Bayes - OoV Solution

> Given a document, classify its author:

$p(\text{Mary Shelly} \mid \text{"The fallen angel becomes a malignant devil"})$

$$P(\text{sentence} \uparrow) P(\text{Mary Shelly})$$

(Naively) Assuming feature independence:

Obtaining an **a-priori** with *Laplace smoothing*:

Counting word occurrences + 1 per class,  
divided by all class words + sentence-words.

$$P(\text{"devil"} \mid \text{Mary Shelly}) = \frac{5 + 1}{1000 + 7}$$

“devil” occurrences + 1  
Total words in Mary Shelly documents + sentence’s tokens

# Classification using Naïve Bayes – Log Probability

> Working with log probability:

$$\hat{c} = \arg \max P(c|d) = \arg \max P(d|c)P(c)$$

For all classes and all words:

$$\operatorname{argmax}_{c \in \text{Classes}} P(c) \prod_{w \in \text{sentence}} P(w_i|c)$$

To ease calculation and speed, we can take logs from both sides:

$$\operatorname{argmax}_{c \in \text{Classes}} \log P(c) + \sum_{w \in \text{sentence}} \log P(w_i|c)$$

# Math Recap: Log rules

$$\log_b(MN) = \log_b(M) + \log_b(N)$$

$$\log_b\left(\frac{M}{N}\right) = \log_b(M) - \log_b(N)$$

$$\log_b(M^p) = p \log_b(M)$$

[Log rules: Justifying the logarithm properties \(article\) | Khan Academy](#)

# Naïve Bayes - Generative

> The bayes formula:

$$\hat{c} = \arg \max P(c|d) = \arg \max \underbrace{P(d|c)}_{\text{Likelihood}} \underbrace{P(c)}_{\text{Prior}}$$

If the prior is known (category/class is given -  $P(c) = 1$ ), we can generate words based on the **Likelihood**:  $P(d|c)$

Reminder:

Discriminative  $\rightarrow P(\textit{class}|\textit{document})$  or

Generative  $\rightarrow P(\textit{document}|\textit{class})$

# Naïve Bayes - Generative

---

## On Discriminative vs. Generative classifiers: A comparison of logistic regression and naive Bayes

---

**Andrew Y. Ng**  
Computer Science Division  
University of California, Berkeley  
Berkeley, CA 94720

**Michael I. Jordan**  
C.S. Div. & Dept. of Stat.  
University of California, Berkeley  
Berkeley, CA 94720

# Naïve Bayes – Final Notes

- > A linear & probabilistic classifier
- > Implementation exists in NLTK & Scikit-Learn (SKL)
- > On SKL one can choose the distribution. The best performance is normally the *Multinomial* Naive Bayes

[A Comparison of Event Models for Naive Bayes Text Classification: binomial.dvi \(washington.edu\)](#)

## A Comparison of Event Models for Naive Bayes Text Classification

Andrew McCallum<sup>‡†</sup>  
mccallum@justresearch.com

<sup>‡</sup>Just Research  
4616 Henry Street  
Pittsburgh, PA 15213

Kamal Nigam<sup>†</sup>  
knigam@cs.cmu.edu

<sup>†</sup>School of Computer Science  
Carnegie Mellon University  
Pittsburgh, PA 15213



# Naïve Bayes – Final Notes

- > Can act as a *discriminative* or *generative* model
- > Still widely used in Linguistics for classification
- > Good summary/deeper dive: [The Optimality of Naive Bayes](#) – sections: Abstract & Naive Bayes and Augmented Naive Bayes (you can safely ignore the augmented Naïve Bayes)

## **The Optimality of Naive Bayes**

**Harry Zhang**

Faculty of Computer Science  
University of New Brunswick

Fredericton, New Brunswick, Canada E3B 5A3  
email: [hzhang@unb.ca](mailto:hzhang@unb.ca)

## Dive deeper

- > [Bayes theorem, the geometry of changing beliefs – YouTube](#)
- > [1.9. Naive Bayes — scikit-learn 1.1.3 documentation](#)
- > [6. Learning to Classify Text \(nltk.org\)](#)
- > [Dan Jurafsky's book – Chapter 4:](#)  
<https://web.stanford.edu/~jurafsky/slp3/4.pdf>