

Learning Distributions

Maximum Likelihood Estimate (MLE)

Bayes Classifier

Aarti Singh

Machine Learning 10-315
Sept 9, 2020



MACHINE LEARNING DEPARTMENT



Logistics

- [Anonymous feedback form](#)
- Recitation on Friday Sept 11 – MLE/MAP + Optimization methods review and hands-on exercises
- QnA1 due TODAY
- HW1 to be released TODAY

Why is ML not ...

➤ Interpolation?

- Noise, stochasticity, transfer across domains, ...

➤ Statistics?

- care about computationally efficiency (feasible, at least polynomial time in input size but typically much faster)

➤ Optimization?

- Don't know true objective function, only stochastic version computed using data samples

$$E_{x,y}[\text{loss}(f(x), y)]$$

➤ Data mining?

- Generalization on new unseen data

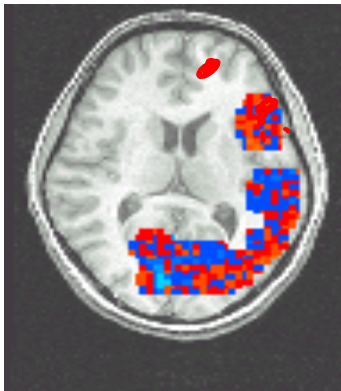
➤ Your question?

Unsupervised Learning

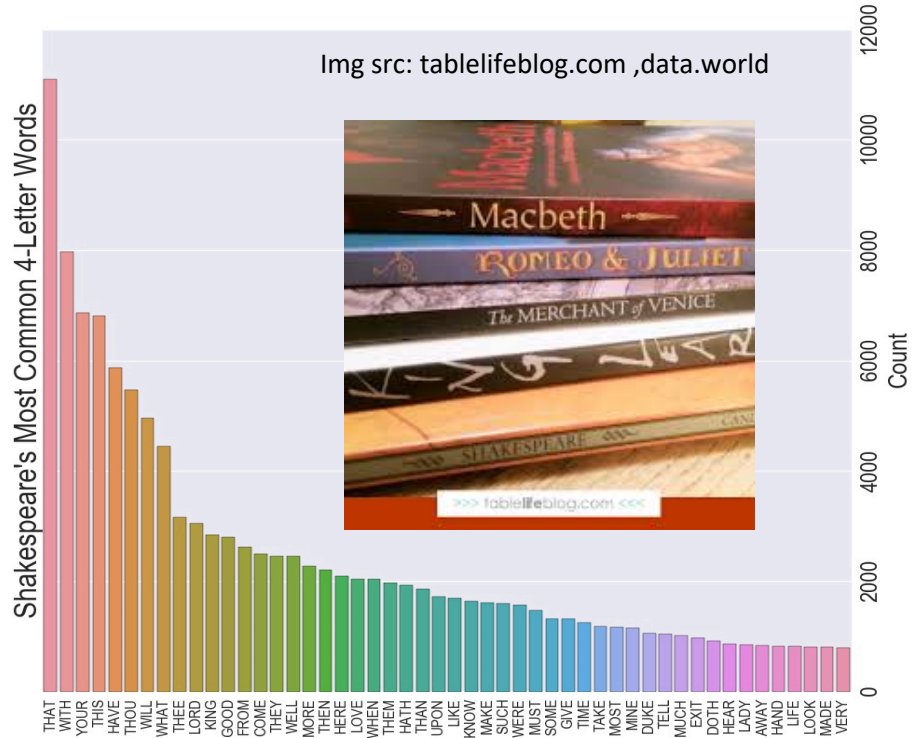
Learning a Distribution



Bias of a coin



voxel

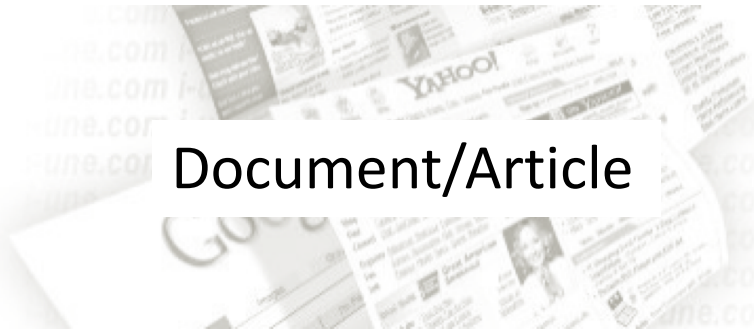


Distribution of words in text

Distribution of brain activity under stimuli

Notion of “Features aka Attributes”

Input $X \in \mathcal{X}$



Document/Article

remember to wake up when class ends

=

wake ends to class remember up when

How to represent inputs mathematically?

- Document vector X ➤ Ideas? $X = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$
 - list of words (different length for each document)
 - frequency of words (length of each document = size of vocabulary), also known as **Bag-of-words** approach ➤ Why might

Misses out context!!

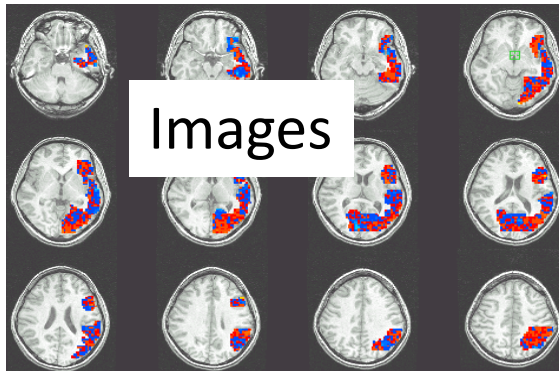
- list of n -grams (n -tuples of words)

$$n=2 \quad d = (\text{Voc})^2 \quad d \uparrow \text{ as } n \uparrow$$

Why might this be limited?

Notion of “Features aka Attributes”

Input $X \in \mathcal{X}$



Input $X \in \mathcal{X}$



How to represent inputs mathematically?

- Image X = intensity/value at each pixel, fourier transform values, SIFT etc.
- Market information X = daily/monthly? price of share for past 10 years

$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} \quad x_i \in \mathbb{R}$$

Distribution of Inputs

Input $X \in \mathcal{X}$

Discrete Probability Distribution $P(X) = P(X=x)$

e.g. $P(\text{head}) = \frac{1}{2}$, $P(\text{word } x \text{ in text}) = p_x$



Probabilities in a distribution sum to 1

$$P(X=x) \geq 0 \quad \sum_x P(X=x) = 1 \quad P(\text{tail}) = 1 - p(\text{head}), \sum_x p_x = 1$$

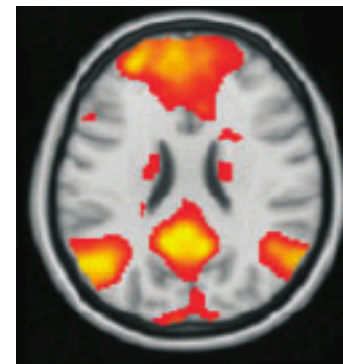
Continuous Probability density $p(x)$

e.g. $p(\text{brain activity})$

$$P(a \leq X \leq b) = \int_a^b p(x) dx$$

Probability density integrate to 1

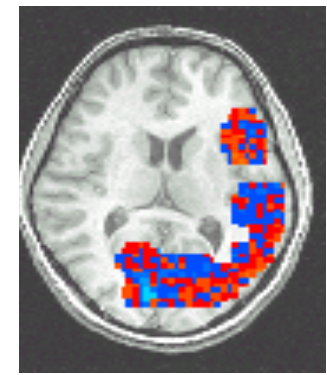
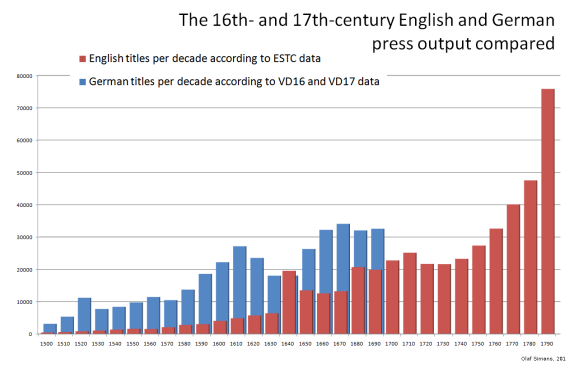
$$p(x) \geq 0 \quad \int p(x) dx = 1$$



Distributions in Supervised tasks

Input $X \in \mathcal{X}$

- Distribution learning also arises in supervised learning tasks e.g. classification
 - ↪ $P(Y = y)$ *topic* Distribution of class labels
 - ↪ $P(X = x | Y = y)$ Distribution of words in 'news' documents
 - ↪ $P(X = x | Y = y)$ Distribution of brain activity under 'stress'

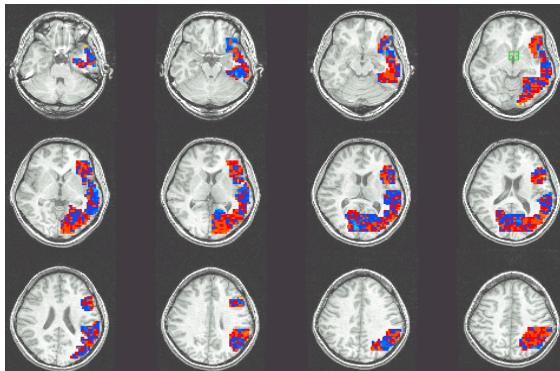


Olaf simons'10

↪ $P(Y = y | X = x)$ Distribution of topics given document

Classification

Goal: Construct prediction rule $f : \mathcal{X} \rightarrow \mathcal{Y}$



High Stress
Moderate Stress
Low Stress

Input feature vector, X

Label, Y

In general: label Y can belong to more than two classes

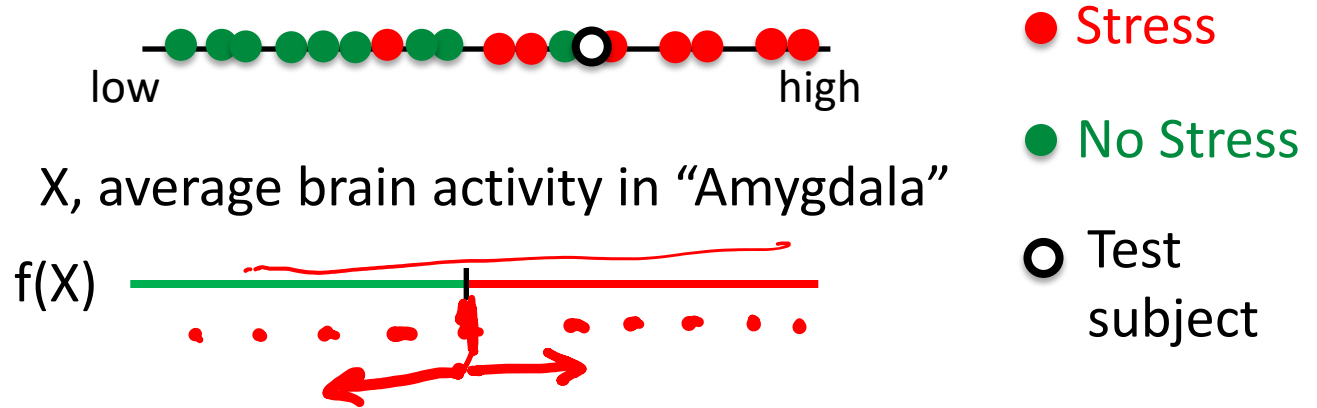
X is multi-dimensional (many features represent an input)

But lets start with a simple case:

label Y is binary (either “Stress” or “No Stress”)

X is average brain activity in the “Amygdala” = $X \in \mathbb{R}$

Binary Classification



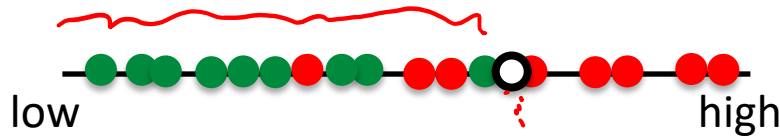
Model X and Y as random variables with joint distribution $P_{XY} = \underline{\underline{P(X, Y)}}$

Training data $\{X_i, Y_i\}_{i=1}^n \sim \underline{\underline{\text{iid}}}$ (independent and identically distributed)
samples from $\underline{\underline{P_{XY}}}$

Test data $\{X, Y\} \sim \underline{\underline{\text{iid}}}$ sample from $\underline{\underline{P_{XY}}}$

Training and test data are independent draws from same distribution

Bayes Classifier

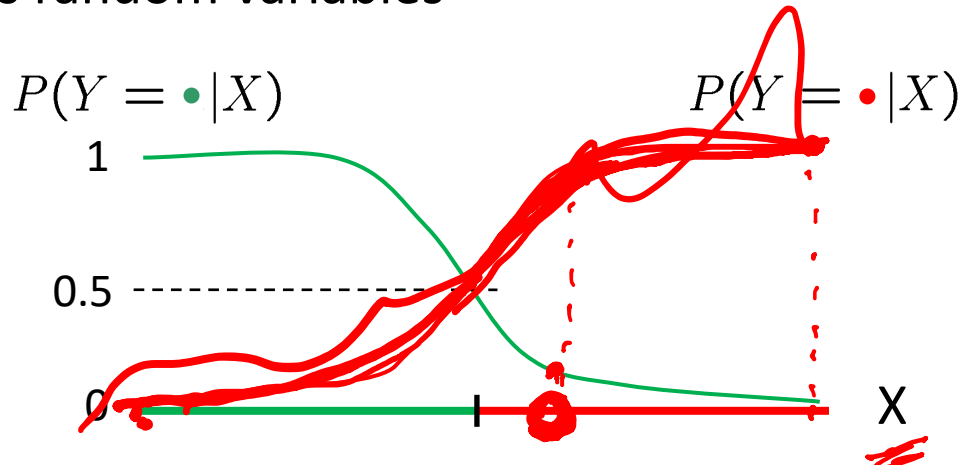


X, average brain activity in "Amygdala"



- Stress
- No Stress
- Test subject

Model X and Y as random variables



For a given X, $f(X) = \text{label } Y \text{ which is more likely}$

$$f(X) = \arg \max_{Y=y} P(Y = y | X = x)$$

$P(Y|X)$

Bayes Rule

Bayes Rule: $P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$

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

To see this, recall:

$$\underline{P(X,Y)} = \underline{P(X|Y)} \underline{P(Y)}$$

$$\underline{P(Y,X)} = \underline{P(Y|X)} \underline{P(X)}$$

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$



Thomas Bayes

Bayes Classifier

Bayes Rule:
$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

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

Bayes classifier:

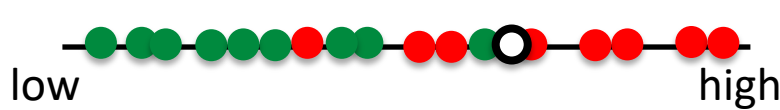
$$f(X) = \arg \max_{Y=y} P(Y = y|X = x) \quad .$$

$$= \arg \max_{Y=y} \underbrace{P(X = x|Y = y)}_{\text{Class conditional}} \underbrace{P(Y = y)}_{\text{Distribution of class}}$$

Class conditional
Distribution of features

Distribution of class

Bayes Classifier



- Stress
- No Stress
- Test subject

X, average brain activity in "Amygdala"



Optimal
E[loss(A(x), y)]
loss = 0/1

$$f(X) = \arg \max_{Y=y} P(X = x | Y = y) P(Y = y)$$

Class conditional

Class distribution

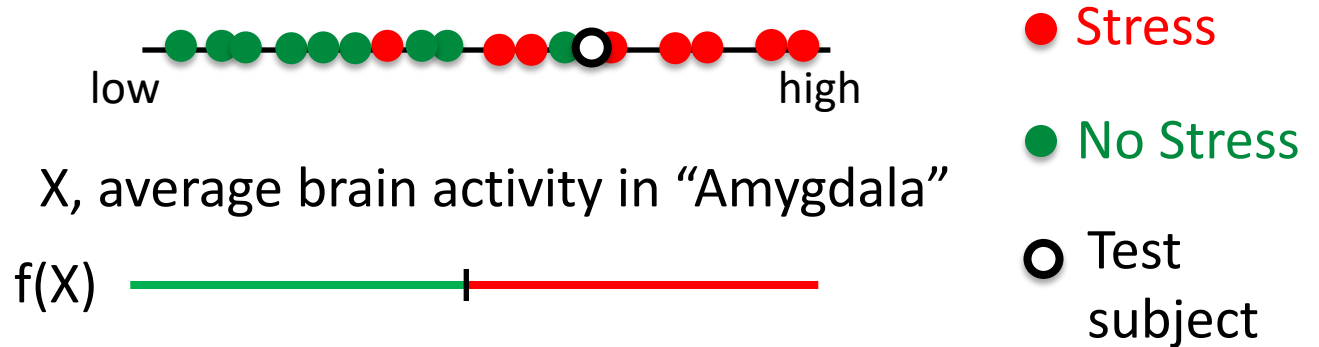
Distribution of features

We can now consider appropriate distribution models for the two terms:

Class distribution $P(Y=y)$ ✓

Class conditional distribution of features $P(X=x | Y=y)$ ✓

Modeling class distribution



Modeling Class distribution $P(Y=y) = \text{Bernoulli}(\theta)$

$$P(Y = \bullet) = \theta \quad \uparrow \text{parameter} \quad P(Y = \bullet) = 1 - \theta$$

Like a coin flip



How to learn parameters from data?

MLE

(Discrete case)

Learning parameters in distributions

$$P(Y = \bullet) = \theta$$

$$P(Y = \bullet) = 1 - \theta$$

Learning θ is equivalent to learning probability of head in coin flip.

➤ How do you learn that?

Data =



Answer: 3/5 ← frequency of heads

➤ Why??

Bernoulli distribution

Data, $D =$



- $P(\text{Heads}) = \theta$, $P(\text{Tails}) = 1 - \theta$
- Flips are **i.i.d.**:
 - **Independent** events
 - **Identically distributed** according to Bernoulli distribution

Choose θ that maximizes the probability of observed data
aka Likelihood

Maximum Likelihood Estimation (MLE)

Choose θ that maximizes the probability of observed data (aka likelihood)

$$\hat{\theta}_{MLE} = \arg \max_{\theta} P(D | \theta)$$

J(θ)
D ~ iid draw from Bel(θ)

MLE of probability of head:

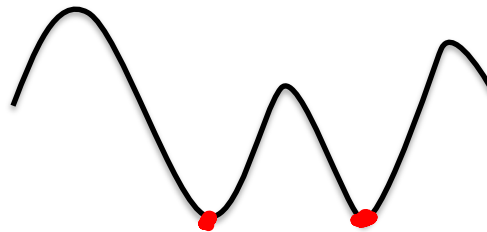
$$\hat{\theta}_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T} = 3/5$$

count of H

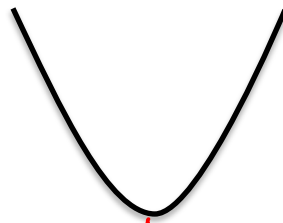
"Frequency of heads"

Short detour - Optimization

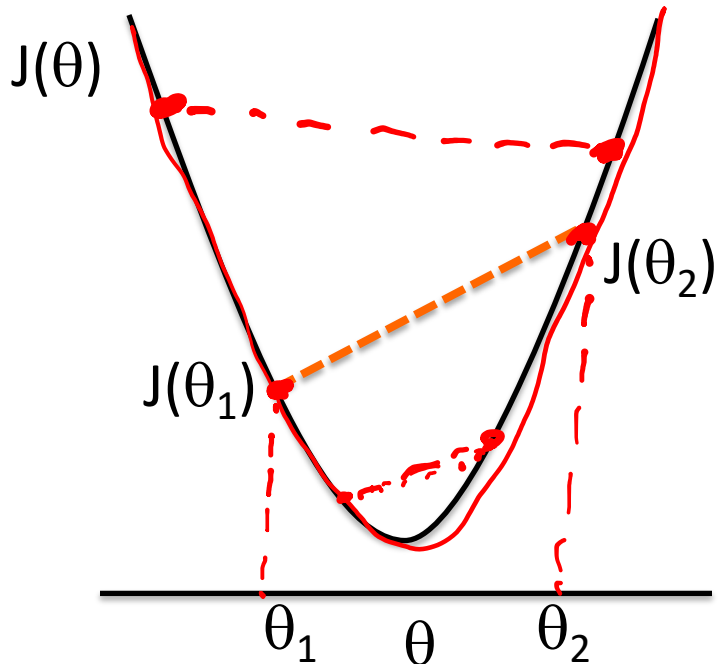
- Optimization objective $J(\theta)$
- Minimum value $J^* = \min_{\theta} J(\theta)$ or $\max_{\theta} J(\theta)$
- Minima (points at which minimum value is achieved) may not be unique



- If function is strictly convex, then minimum is unique

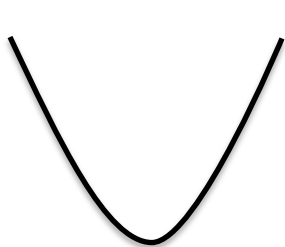


Convex functions

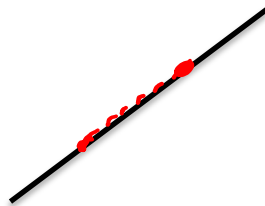


A function $J(\theta)$ is called **convex** if the line joining two points $J(\theta_1), J(\theta_2)$ on the function does not go below the function on the interval $[\theta_1, \theta_2]$

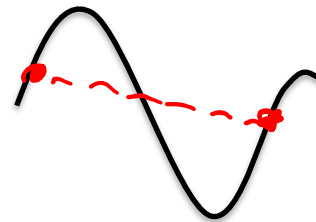
(Strictly) Convex functions
have a unique minimum!



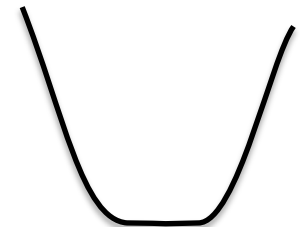
Convex



Both Concave
& Convex



Neither

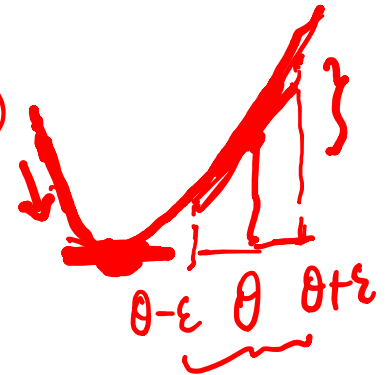


Convex but not
strictly convex²¹

Optimizing convex (concave) functions

- Derivative of a function

$$\frac{\Delta J(\theta)}{\Delta \theta} = \lim_{\epsilon \rightarrow 0} \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$$



- Partial derivative

$$\frac{\partial J(\theta)}{\partial \theta}$$

- Derivative is zero at minimum of a convex function

$$g(\theta) = \left. \frac{\partial J(\theta)}{\partial \theta} \right|_{\theta = \theta^*} = 0$$

- Second derivative is positive at minimum of a convex function

$$\frac{\partial g(\theta)}{\partial \theta} = \frac{\partial^2 J(\theta)}{\partial \theta^2} \geq 0$$

Optimizing convex (concave) functions

➤ What about

concave functions?

non-convex/non-concave functions?

functions that are not differentiable?

optimizing a function over a bounded domain aka
constrained optimization?



Maximum Likelihood Estimation (MLE)

Choose θ that maximizes the probability of observed data (aka likelihood)

$$\hat{\theta}_{MLE} = \arg \max_{\theta} \underbrace{P(D | \theta)}_{\mathcal{J}(\theta)}$$

MLE of probability of head:

① $D \sim \text{iid}$
② $\text{Ber}(\theta)$

$$\hat{\theta}_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T} = 3/5$$

"Frequency of heads"