Announcements

- Anonymous feedback form
- Recitation on Friday Sept 10 Convexity review
- QnA1 due TODAY
- HW1 to be released TODAY

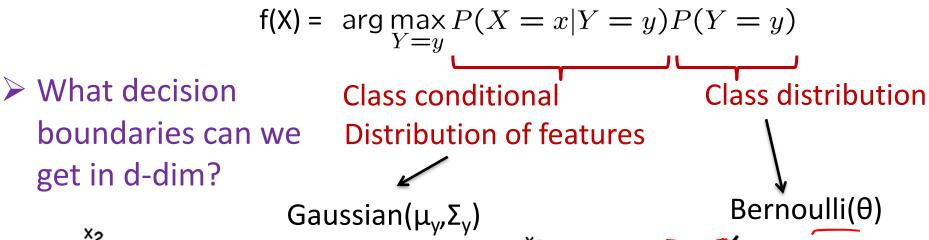
Recap

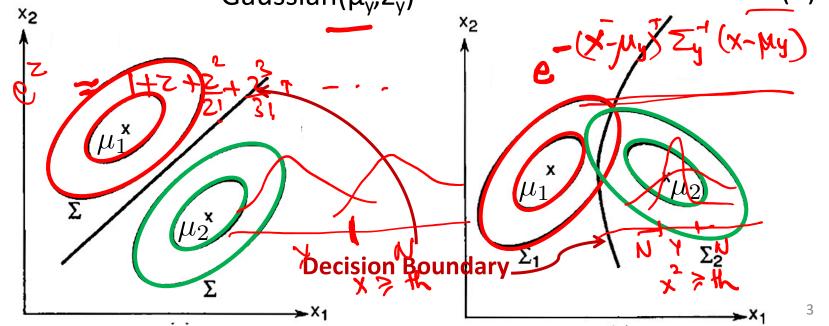
Bayes classifier – assumes P_{XY} known, optimal for 0/1 loss

$$f(X) = \arg\max_{Y=y} P(Y=y|X=x) = \mathop{\rm E}_{xy} [1+(x) \neq y] \\ = \arg\max_{Y=y} P(X=x|Y=y) P(Y=y) \\ {\rm Class \ conditional} {\rm Class \ distribution} \\ {\rm Distribution \ of \ features}$$

- Gaussian Bayes classifier assumes
 Class distribution is Bernoulli/Multinomial
 Class conditional distribution of features is Gaussian
- Decision boundary (binary classification)

d-dim Gaussian Bayes classifier





Decision Boundary of Gaussian Bayes

Decision boundary is set of points x: P(Y=1|X=x) = P(Y=0|X=x)

Compute the ratio

Compute the ratio
$$1 = \frac{P(Y = 1 | X = x)}{P(Y = 0 | X = x)} = \frac{P(X = x | Y = 1)P(Y = 1)}{P(X = x | Y = 0)P(Y = 0)} = \frac{P(X = x | Y = 0)P(Y = 0)}{P(X = x | X = 0)} = \frac{P(X = x | Y = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | Y = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | Y = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | Y = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | Y = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | Y = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | Y = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | Y = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | Y = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | Y = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | Y = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | Y = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | Y = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | Y = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | Y = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | Y = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | X = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | X = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | X = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | X = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | X = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | X = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | X = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | X = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | X = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | X = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | X = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | X = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | X = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | X = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | X = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | X = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | X = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | X = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | X = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | X = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | X = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | X = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X = x | X = 0)P(X = 0)}{P(X = x | X = 0)} = \frac{P(X$$



$$= \sqrt{\frac{|\Sigma_0|}{|\Sigma_1|}} \exp\left(-\frac{(x-\mu_1)^{\mathsf{T}}\Sigma_1^{-1}(x-\mu_1)}{2} + \frac{(x-\mu_0)^{\mathsf{T}}\Sigma_0^{-1}(x-\mu_0)}{2}\right) \underbrace{\frac{\theta}{1-\theta}}_{1-\theta}$$

In general, this implies a quadratic equation in x. But if $\Sigma_1 = \Sigma_0$, then quadratic part cancels out and decision boundary is linear.

Glossary of Machine Learning

- Feature/Attribute
- iid
- Bayes classifier
- Class distribution
- Class conditional distribution of features
- Decision boundary

Aarti Singh

Machine Learning 10-315 Sept 8, 2021





How to learn parameters from data? MLE

(Discrete case)

Learning parameters in distributions

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

$$P(Y = 0) = 1 - \theta$$

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

➤ How do you learn that?

Answer: 3/5

➤ Why??

Bernoulli distribution

- Parameter θ : P(Heads) = θ , P(Tails) = 1- θ
- Flips are i.i.d.:
 - Independent events
 - Identically distributed according to Bernoulli distribution

<u>Choose θ that maximizes the probability of observed data</u> <u>aka Likelihood</u>

Choose θ that maximizes the probability of observed data (aka likelihood) $D = X_1 \dots X_n$ $X_i = \begin{cases} 0 & T \end{cases}$

$$\widehat{\theta}_{MLE} = \arg\max_{\theta} P(D \mid \theta)$$

MLE of probability of head:

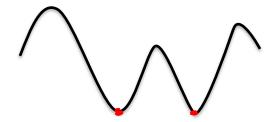
bability of head:
$$\widehat{\theta}_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T} = 3/5$$
 "Frequency of heads"

"Frequency of heads"

 $= P(X_1 ... X_n; \theta)$

Short detour - Optimization

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

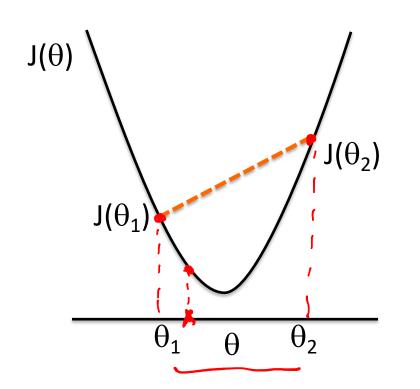


• If function is strictly convex, then minimum is unique



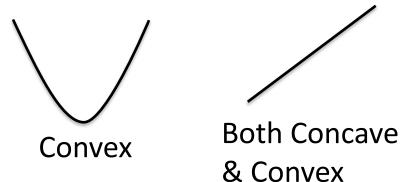
Convex functions

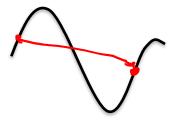




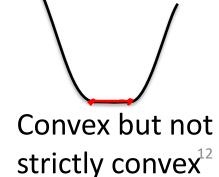
A function $J(\theta)$ is called **convex** if the line joining two points $J(\theta_1), J(\theta_2)$ on the function does not go below the function on the interval $[\theta_1, \theta_2]$

(Strictly) Convex functions have a unique minimum!





Neither



Optimizing convex (concave) functions

Derivative of a function

$$\frac{\partial J(0)}{\partial l(0)} = \lim_{k \to 0} \frac{J(0+k) - J(0)}{k}$$

0 Oth -ve - 0 - tre

- Partial derivative
- Derivative is zero at minimum of a convex function

Second derivative is positive at minimum of a convex function

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?

Bernoulli MLE Derivation

$$\widehat{\theta}_{MLE} = \arg\max_{\theta} P(D \mid \theta)$$

$$= \arg\max_{\theta} \widehat{\prod_{i=1}^{N} \theta^{X_i}} (1-\theta)^{1-X_i}$$

$$= \arg\max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \max_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \min_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \min_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \min_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \min_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \min_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \min_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \min_{\theta} \widehat{\prod_{i=1}^{N} Y_i \log \theta} + (1-X_i) \log (1-\theta)$$

$$= 2\pi \min_{\theta} \widehat{\prod_{i=1}^{N$$

Multinomial distribution

Data, D = rolls of a dice



- $P(1) = p_1$, $P(2) = p_2$, ..., $P(6) = p_6$ $p_1 + + p_6 = 1$
- Rolls are i.i.d.:
 - Independent events
 - Identically distributed according to Multinomial(θ) distribution where

$$\theta = \{p_1, p_2, ..., p_6\}$$

<u>Choose θ that maximizes the probability of observed data</u> <u>aka "Likelihood"</u>

Choose θ that maximizes the probability of observed data

$$\widehat{\theta}_{MLE} = \arg \max_{\theta} P(D \mid \theta)$$

MLE of probability of rolls:

$$\hat{\theta}_{MLE} = \hat{p}_{1,MLE}, \dots, \hat{p}_{6,MLE}$$

$$\hat{p}_{y,MLE} = \frac{\alpha_y - \text{Rolls that turn up y}}{\sum_y \alpha_y}$$
 Total number of rolls

"Frequency of roll y"

How to learn parameters from data? MLE

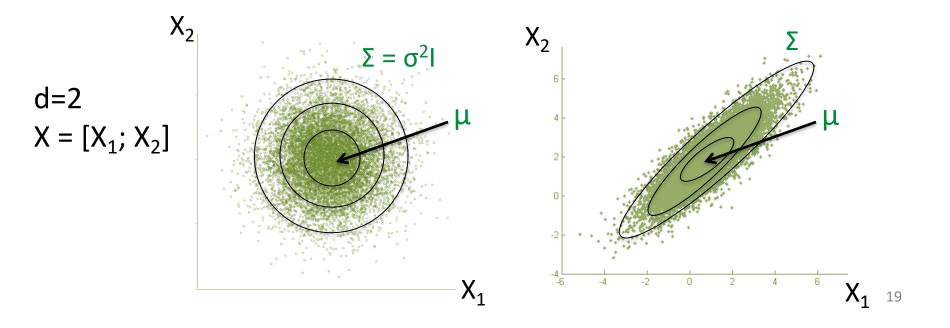
(Continuous case)

d-dim Gaussian distribution

X is Gaussian $N(\mu, \Sigma)$

 μ is d-dim vector, Σ is dxd dim matrix

$$P(X = x | \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right),$$



How to learn parameters from data? MLE

(Continuous case)

Gaussian distribution

- Parameters: μ mean, σ^2 variance
- Data are i.i.d.:
 - Independent events
 - Identically distributed according to Gaussian distribution

Choose θ = (μ , σ ²) that maximizes the probability of observed data

$$\begin{split} \widehat{\theta}_{MLE} &= \arg\max_{\theta} \ P(D \mid \theta) \\ &= \arg\max_{\theta} \prod_{i=1}^{n} P(X_i | \theta) \quad \text{Independent draws} \\ &\sum_{i=1}^{n} P(X_i | \theta) \quad \text{Independent$$

Choose θ = (μ , σ ²) that maximizes the probability of observed data

$$egin{array}{ll} \widehat{ heta}_{MLE} &= \arg\max_{ heta} \; P(D \mid heta) \\ &= \arg\max_{ heta} \prod_{i=1}^n P(X_i | heta) \quad & \text{Independent draws} \\ &= \arg\max_{ heta} \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(X_i - \mu)^2/2\sigma^2} \quad & \text{Identically distributed} \end{array}$$

Choose $\theta = (\mu, \sigma^2)$ that maximizes the probability of observed data

$$\begin{split} \widehat{\theta}_{MLE} &= \arg\max_{\theta} \ P(D \mid \theta) \\ &= \arg\max_{\theta} \prod_{i=1}^n P(X_i | \theta) \quad \text{Independent draws} \\ &= \arg\max_{\theta} \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(X_i - \mu)^2/2\sigma^2} \quad \text{Identically distributed} \\ &= \arg\max_{\theta = (\mu, \sigma^2)} \frac{1}{(2\pi\sigma^2)^{n/2}} e^{-\sum_{i=1}^n (X_i - \mu)^2/2\sigma^2} \quad \text{and} \quad J(\theta) \end{split}$$

MLE for Gaussian mean

> Poll

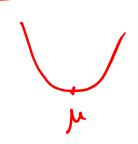
$$P(D|\theta) = \frac{1}{(2\pi\sigma^2)^{n/2}} e^{-\sum_{i=1}^{n} (X_i - \mu)^2 / 2\sigma^2}$$

A.
$$\max_{\mu} \sum_{i=1}^{n} (X_i - \mu)^2$$

c.
$$\max_{\mu} \mu^2 - 2\mu \sum_{i=1}^{n} X_i$$

B.
$$\min_{\mu} \sum_{i=1}^{n} (X_i - \mu)^2$$

D.
$$\max_{\mu} n\mu^2 - 2\mu \sum_{i=1}^{n} X_i$$



MLE for Gaussian mean and variance

$$\widehat{\mu}_{MLE} = \frac{1}{n} \sum_{i=1}^{n} x_i \qquad \checkmark$$

$$\widehat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \widehat{\mu})^2$$

Self exercise:

Derive MLE of variance?

d-dimensional versions?



More coming up in HW1

Unbiased estimates

- Bias of an estimate $\mathsf{E}[\widehat{ heta}]$ heta
- Unbiased estimate if $E[\hat{\theta}] = \theta$
 - ➤ Is the MLE of mean unbiased?

➤ Is the MLE of variance unbiased?

$$E[\hat{\sigma}_{\text{MLE}}] = \hat{\sigma}_{n}^{2}$$

How can you make it unbiased?