# INTRODUCTION TO MACHINE LEARNING COMPSCI 4ML3

## LECTURE 9

### HASSAN ASHTIANI

# REVIEW: BAYES RULE, CHAIN RULE

$P(X)$

$P(X = a)$

- Joint distribution: $P(X, Y)$

- Sum rule: $P(X) = \sum_Y P(X, Y)$

- Conditional distribution: $P(X|Y) = \dfrac{P(X, Y)}{P(Y)}$

- Bayes rule: $P(X|Y) = \dfrac{P(Y|X)P(X)}{P(Y)}$

- Chain Rule: $P(X_1, X_2, \ldots, X_k) = $

  $P(X_1) P(X_2|X_1) P(X_3|X_2, X_1), \ldots, P(X_k|X_1, X_2, \ldots X_{k-1})$

- $P(X_1, X_2, \ldots, X_k|Y) = P(X_1|Y) P(X_2|X_1, Y) \ldots$

- 

  $P(X|Y), P(Y|X) \to P(X, Y) = ?$

# REVIEW: INDEPENDENCE

$$P(X|Y) = P(X)$$

- $X$ AND $Y$ ARE INDEPENDENT IF $P(X, Y)$ $= P(X)P(Y)$

ASSUME $X_1, \ldots, X_k$ ARE INDEPENDENT GIVEN Y

- $P(X_1, X_2, \ldots, X_k | Y) =$

$$P(X_1|Y) \, P(X_2|Y) \cdots P(X_k|Y)$$

-

# STATISTICAL APPROACH TO ML

- Our goal is to do well on new/unseen (test) data

- We were mostly minimizing the training error

  - Directly/systematically optimizing the test error?

- There is uncertainly about the unseen data

  - We cannot be 100% sure about the performance of <u>any method</u> on the test data

  - A method that works well on test set **<u>most of the time</u>**?

# I.I.D ASSUMPTION

- Assume there is an underlying (unknown) distribution $D$

- Assume that each of the training and test instances are sampled independently from $D$

- We say train and test sets are I.I.D. (independent and identically distributed) samples generated from $D$

# I.I.D ASSUMPTION

- Why are these assumptions necessary?
  - Same distribution for all samples ("identically")
  - Independent samples ("Independently")
  - Same distribution for train and test
  - The distribution is unknown

# PARAMETER ESTIMATION

- ASSUME THAT THE DISTRIBUTION COMES FROM SOME KNOWN FAMILY

  - BERNOULLI, GAUSSIAN, …

- USE THE TRAINING SET TO ESTIMATE THE VALUE OF THE UNKNOWN DISTRIBUTION PARAMETERS

- USEFUL IN BOTH UNSUPERVISED AND SUPERVISED LEARNING

# BIASED COIN EXAMPLE (UNSUPERVISED)

- FLIPPING A COIN
  - OUTCOME IS HEAD $(0)$ OR TAIL $(1)$, SO $x \in \{0,1\}$
- $P(x = 0) = \alpha, P(x = 1) = 1 - \alpha$
  - $x$ IS A BERNOULLI RANDOM VARIABLE
- BIAS $(\alpha)$ IS UNKNOWN (THE PARAMETER)
- GIVEN AN I.I.D SAMPLE, **ESTIMATE** $\alpha$
  - $X = (x^1, x^2, x^3, \ldots, x^n)$
  - E.G., $n = 10, X = (0,0,1,1,0,1,0,1,0,0)$

$\hat{\alpha} = \dfrac{6}{10}$

# ESTIMATING THE BIAS OF THE COIN

- LET $n_0 = \#$HEADS, $n_1 = \#$TAILS (SO $n_0 + n_1 = n$)

- IS $\hat{\alpha} = \frac{n_0}{n_0 + n_1}$ A GOOD ESTIMATE?

- IS THERE A RATIONAL BEHIND THIS ESTIMATE?

# MAXIMUM LIKELIHOOD ESTIMATE (MLE)

GIVEN THE TRAINING SET $X = (x^1, x^2, \ldots, x^n)$, ESTIMATE $\alpha$.

- MLE MAXIMIZES THE PROBABILITY OF THE OBSERVATIONS GIVEN THE PARAMETERS

  $\rightarrow$ likelihood

  - $$\alpha^{ML} = \underset{\alpha}{argmax}\, P\,(X|\alpha)$$

- EQUIVALENTLY (WHY?)

  - $$\alpha^{ML} = \underset{\alpha}{argmin} -\log\left(\sum_i P\,(x^i|\alpha)\right)$$

  $$-\sum \log\left(P(x^i|\alpha)\right)$$

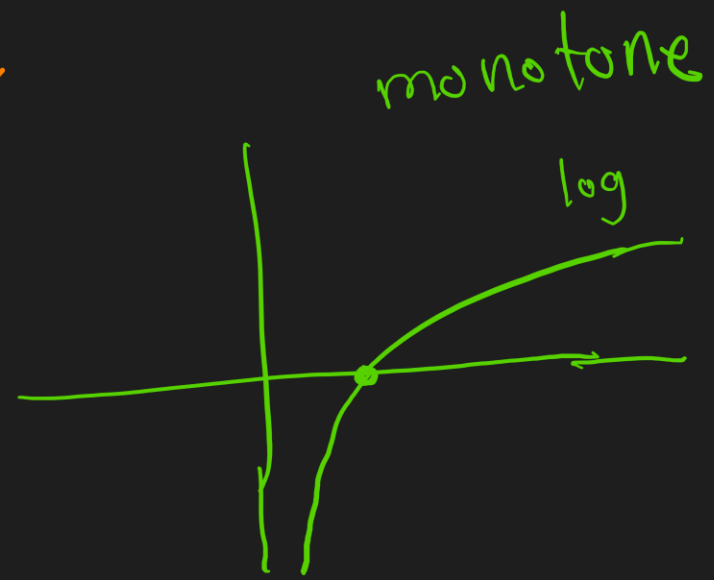# NEGATIVE-LOG-LIKELIHOOD

$$\alpha_{ML} = \arg\max_{\alpha} P(X|\alpha) = \arg\min_{\alpha} - P(X|\alpha)$$
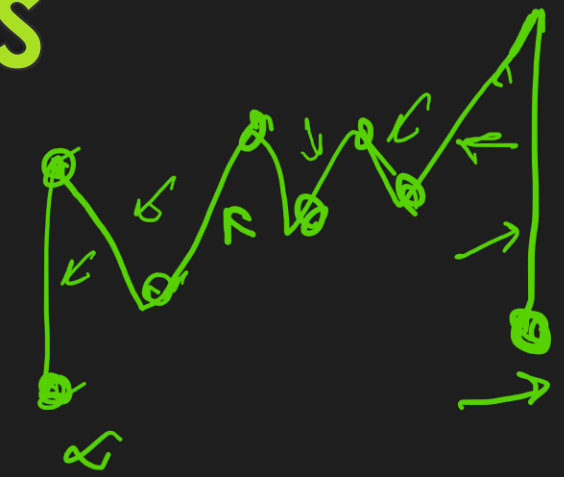
$$= \arg\min_{\alpha} -\log P(X|\alpha)$$

$$= \arg\min_{\alpha} -\log \prod_{i=1}^{n} P(x^i|\alpha)$$

$$= \arg\min_{\alpha} -\sum_{i=1}^{n} \log P(x^i|\alpha) \checkmark$$

monotone

log

# MAXIMUM LIKELIHOOD FOR COINS

$$\text{Likelihood} = p(X|\alpha) = \prod_{i=1}^{n} p(x^i|\alpha) =$$

$$= \left( \prod_{i:x^i=0} p(x^i|\alpha) \right) \left( \prod_{i:x^i=1} \frac{p(x^i|\alpha)}{1-\alpha} \right)$$

$$= \alpha^{n_0} (1-\alpha)^{n_1} = f(\alpha)$$

$$0 = \frac{\partial f}{\partial \alpha} = \quad n_0 \alpha^{n_0-1} (1-\alpha)^{n_1} - \alpha^{n_0} n_1 (1-\alpha)^{n_1-1}$$

$$\Rightarrow n_0 (1-\alpha) - n_1 \alpha = 0 \rightarrow \boxed{\alpha^{ML} = \frac{n_0}{n_0 + n_1}}$$

# MAXIMUM A POSTERIORI ESTIMATE

- Maximizes the probability of the parameters given the observations

$$\alpha^{MAP} = \underset{\alpha}{argmax}\, P\,(\alpha|X)$$

- $$\alpha^{MAP} = \underset{\alpha}{argmin}\left(-\log(P(\alpha)) - \sum_{i=1}^{n} \log P(x^i|\alpha)\right)$$

# PRIOR VS POSTERIOR DISTRIBUTIONS

- $P(\alpha)$ CAPTURES THE **PRIOR** DISTRIBUTION

- $P(\alpha|X)$ CAPTURES THE **POSTERIOR** DISTRIBUTION

- IN OTHER WORDS,

  - WE START BY A PRIOR BELIEF ABOUT VALUE OF $\alpha$

  - OUR BELIEF IS UPDATED AFTER SEEING SOME REAL DATA

  - THIS IS A **BAYESIAN** APPROACH

# MAP FOR COINS – UNIFORM PRIOR

# MAP FOR COINS – NONUNIFORM PRIOR