

ECE 493: Reinforcement Learning

Probability and Stats Review

Mark Crowley

Spring 2020

1 Probability and Statistics Review

- Probability Definitions
- Bayes Theorem
- Entropy
- Probabilistic Distance Metrics

Probability and Statistics Review

- Factoring Probability Distributions
- The joint prob decomposes into multiplication of probs if vars are indep
- Entropy
- KL-divergence, Mahalanobis distance

Joint and Conditional Probability

Given event X (binary or multivalued). $p(X = x) = p(x)$ is the probability of the event that X takes on the value x .

Probability of A or B occurring:

$$\begin{aligned} p(A \vee B) &= p(A) + p(B) - p(A \wedge B) \\ &= p(A) + p(B) \text{ if } A \text{ and } B \text{ are mutually exclusive} \end{aligned}$$

Joint Probabilities: Product Rule and Chain Rule

$$\begin{aligned} p(A, B) &= p(A \wedge B) = p(A|B)p(B) = p(B|A)p(A) \\ p(X_1, X_2, \dots, X_D) &= p(X_1)p(X_2|X_1)p(X_3|X_2, X_1) \dots p(X_D|X_{1:D-1}) \end{aligned}$$

Marginal and Conditional Probability

Marginal Distribution:

$$p(A) = \sum_b p(A, B) = \sum_b p(A|B = b)p(B = b)$$

Conditional Probability:

$$p(A|B) = \frac{p(A, B)}{p(B)} \text{ if } p(B) > 0$$

“Probability of A given B ”

Bayes Theorem

Given a hypothesis h and observed evidence e :

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$

$$p(h|e) = \frac{p(e|h)p(h)}{p(e)}$$

$$p(\text{cancer}|\text{testresult}) = \frac{p(\text{testresult}|\text{cancer})p(\text{cancer})}{p(\text{testresult})}$$

- *An aside:* Bayesian Statistics vs Frequentist Statistics
- very important, for knowing how to update a model based on new evidence, also tells you how to turn around a $p(X|Y)$ into a $p(Y|X)$

Bayes Theorem For Multidimensional Data

For data x with prediction target y :

$$\begin{aligned}\text{posterior} &= \frac{\text{likelihood} \times \text{prior}}{\text{evidence}} \\ p(y|x_1, \dots, x_n) &= \frac{p(x_1, \dots, x_n|y)p(y)}{p(x_1, \dots, x_n)}\end{aligned}$$

If we knew that all of the features x_i were **independent** then we'd have:

$$p(y|x_1, \dots, x_n) = \frac{p(y) \prod_{i=1}^n p(x_i|y)}{p(x_1, \dots, x_n)}$$

Bayes Theorem as a Proportion

The probability of the evidence is constant and just for normalizing to a probability. So if we only want to compare probabilities we can drop it:

$$p(y|x_1, \dots, x_n) = \frac{p(y) \prod_{i=1}^n p(x_i|y)}{p(x_1, \dots, x_n)}$$

$$p(y|x_1, \dots, x_n) \propto p(y) \prod_{i=1}^n p(x_i|y)$$

Naive Bayes Classification

$$p(y|x_1, \dots, x_n) \propto p(y) \prod_{i=1}^n p(x_i|y)$$

The **Naive Bayes classifier** uses this form to estimate the label y making the “naive” independence assumption.

$$\hat{y} = \arg \max_y p(y) \prod_{i=1}^n p(x_i|y)$$

- $p(y)$ can be estimated simply with counts of the frequency of each class in the data
- $p(x_i|y)$ is a known distribution you specify: Gaussian, Multinomial, Bernoulli

Unconditional Independence

If two random variables variable X and Y are independent we denote it as $X \perp Y$

$$X \perp Y \text{ iff } p(X, Y) = p(X)p(Y)$$

Outline of

1 Probability and Statistics Review

- Probability Definitions
- Bayes Theorem
- **Entropy**
- Probabilistic Distance Metrics

$$H(X) = - \sum_{x \in X} p(x) \log_2 p(x)$$

- The higher the entropy the higher the uncertainty for that value.
- Also measures surprise of seeing the observation.
- How much information is represented by this observation.

Visualizing Entropy

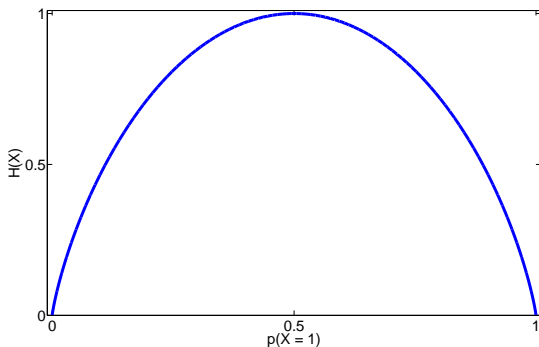


Figure: Binary Entropy Function: Entropy of the Bernoulli random variable as a function of θ . The maximum entropy is $\log_2 2 = 1$ when $\theta = 0.5$ (i.e. when the distribution is uniform).

KL-Divergence or Relative Entropy

Kullback-Leibler Divergence (KL-Divergence) is a common method for measuring the dissimilarity between two probability distributions P and Q .

$$KL(P||Q) = \sum_{i=1}^N P(i) \log \frac{P(i)}{Q(i)}$$

- $KL(P||Q) \geq 0$ and equals zero iff $P = Q$
- How much information you'd lose approximating Q with P
- In general $KL(P||Q) \neq KL(Q||P)$

Outline of

1 Probability and Statistics Review

- Probability Definitions
- Bayes Theorem
- Entropy
- Probabilistic Distance Metrics

Mahalanobis Distance

Another way to measure difference between vectors that accounts for their distribution

$$d(x, y) = \sqrt{(x - y)^T S^{-1} (x - y)}$$

Where x and y share the *same* distribution and covariance matrix S

Interpretation: Multi-dimensional generalization of measuring how many standard deviations away X is from the mean of Y . Disimilarity between two vectors.

- Distance is preserved under linear transformations of data.
- Distance is zero if (x_i, y_i) is at the mean of D , and grows as it moves away from the mean.
- If $S = I$ then equivalent to Euclidean distance.

Mutual Information (MI)

The **mutual information (MI)** between two vectors X, Y measures how similar the joint distribution $p(X, Y)$ is to the factored distribution $p(X)p(Y)$:

$$MI(X, Y) = \sum_{x \in X} \sum_{y \in Y} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}$$

- $MI(X, Y)$ is always nonnegative
- Equals 0 iff X, Y are independent
- Notice this is just the KL-Divergence between the distributions $p(X, Y)$ and $p(X)p(Y)$

[From [?]]

Relation of MI to Entropy

The entropy $H(X)$ and mutual information are related:

$$H(X) = - \sum_{x \in X} p(x) \log p(x) \quad (1)$$

$$MI(X, Y) = H(X) + H(Y) - H(X, Y) \quad (2)$$

- MI can be seen as the *reduction in entropy* on the labels that results from observing feature value x_j
- Some measures use MI normalized by the entropy $H(X)$

[From [?]]

Another measure you could use is **information gain**.

$$IG(Y, X) = H(Y) - H(Y|X)$$

Outline of

1 Probability and Statistics Review

- Probability Definitions
- Bayes Theorem
- Entropy
- Probabilistic Distance Metrics

Hypothesis Testing

- Given a known distribution D_0 we think produced the data, call this our **null hypothesis** (often denoted H_0)
- Want to ask whether we can *reject the null hypothesis* given some observed data.
- Say D_0 is $N(0, 1)$ a standardized Gaussian and the sample is $x = 2.576$.
- $p(|u| \leq 2.576) = .99$: The probability of a sample u taken from $N(0, 1)$ being less than 2.576 is 99%.
- So we say the difference of the sample x from the assumed distribution is *statistically significant*
- Also, say that the sample x lets us “reject the null hypothesis at the 0.1 confidence level”.
- Many methods for doing this, for discrete data one is the Chi-squared (χ^2) Test.

Chi-squared (χ^2) Test

χ^2 statistics can be used to test whether a feature is statistically significant in predicting a class. For a feature x_f and class y_k we can formulate the Chi-square test

$$\begin{aligned}\chi^2(x_{fi}, y_k) &= \sum_{x_{fi} \in X_f} \sum_{y_k \in Y} \frac{(O_{ik} - E_{ik})^2}{E_{ik}} \\ &= N \sum_{x_{fi} \in X_f} \sum_{y_k \in Y} p_i p_k \left(\frac{(O_{ik}/N) - p_i p_k}{p_i p_k} \right)^2\end{aligned}$$

- O are the observed counts of joint events and E are their expected counts.
- χ^2 tests the hypothesis that the features and the classes are assigned randomly and independent
- The higher the value of χ^2 , the more likely we reject the null hypothesis of independent, random assignment of classes.
- Thus, the higher the value of χ^2 the more likely this feature f gives a statistically significant discrimination between the classes.

Contingency Table

The counts for χ^2 can be obtained using a contingency table

		y_k	
		1	0
x_{fi}	1	o_{11}	o_{12}
	0	o_{21}	o_{22}

- o_{11} is number of samples in the class that has the feature
- o_{21} is number of samples in the class that doesn't have the feature
- o_{12} is number of samples in other classes that has the feature
- o_{22} is number of samples in other classes that doesn't have the feature

Contingency Table

		y_k				E_{ik}	
		1	0			1	0
x_{fi}	1	100	70	= 170	1	104.6	65.4
	0	60	30	= 90	0	55.4	34.6
		= 160	= 100	N = 260			

$$E_{11} = (o_{11} + o_{21})(o_{11} + o_{12})/N \quad E_{12} = (o_{12} + o_{22})(o_{11} + o_{12})/N \quad (3)$$

$$E_{21} = (o_{11} + o_{21})(o_{21} + o_{22})/N \quad E_{22} = (o_{12} + o_{22})(o_{21} + o_{22})/N \quad (4)$$

Contingency Table

		y_k				E_{ik}		
		1	0			1	0	
x_{fi}	1	100	70	= 170	1	104.6	65.4	
	0	60	30	= 90	0	55.4	34.6	
		= 160				= 100		N = 260

$$\chi^2 = \frac{(100 - 104.6)^2}{104.6} + \frac{(70 - 65.4)^2}{65.4} + \frac{(60 - 55.4)^2}{55.4} + \frac{(30 - 34.6)^2}{34.6} = 1.51 \quad (5)$$

The number of *degrees of freedom* here is 1. Now we can use a Chi-squared lookup table to find the critical value for this number at a desired significance level. We see that for $p=0.05$ we need $\chi^2 > 3.8$ to reject the null hypothesis and claim that our feature is significant. In this case we can only claim $p=0.30$ significance.

χ^2 Lookup Table

Degrees of freedom (df)	χ^2 value ^[18]										
1	0.004	0.02	0.06	0.15	0.46	1.07	1.64	2.71	3.84	6.64	10.83
2	0.10	0.21	0.45	0.71	1.39	2.41	3.22	4.60	5.99	9.21	13.82
3	0.35	0.58	1.01	1.42	2.37	3.66	4.64	6.25	7.82	11.34	16.27
4	0.71	1.06	1.65	2.20	3.36	4.88	5.99	7.78	9.49	13.28	18.47
5	1.14	1.61	2.34	3.00	4.35	6.06	7.29	9.24	11.07	15.09	20.52
6	1.63	2.20	3.07	3.83	5.35	7.23	8.56	10.64	12.59	16.81	22.46
7	2.17	2.83	3.82	4.67	6.35	8.38	9.80	12.02	14.07	18.48	24.32
8	2.73	3.49	4.59	5.53	7.34	9.52	11.03	13.36	15.51	20.09	26.12
9	3.32	4.17	5.38	6.39	8.34	10.66	12.24	14.68	16.92	21.67	27.88
10	3.94	4.87	6.18	7.27	9.34	11.78	13.44	15.99	18.31	23.21	29.59
P value (Probability)	0.95	0.90	0.80	0.70	0.50	0.30	0.20	0.10	0.05	0.01	0.001

Figure: From wikipedia:Chi-squared Distribution