



ISTA 421/521

Introduction to Machine Learning

Lecture 6: Probability Review, Expectation, Discrete Prob. Distributions

Clay Morrison

clayton@sista.arizona.edu

Gould-Simpson 819

Phone 621-6609

11 September 2014

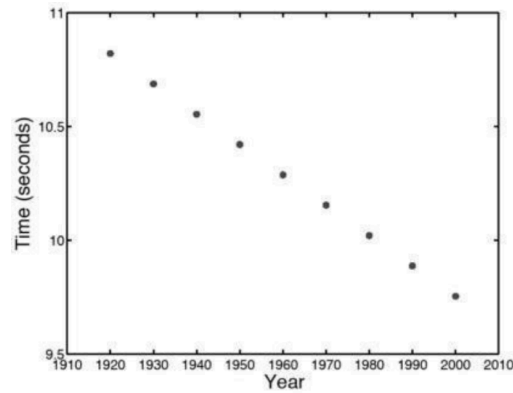
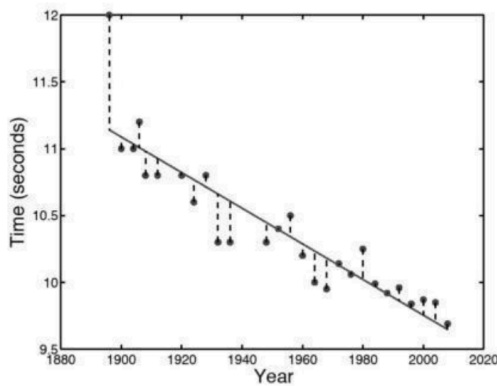


Next Topics

- Probability Basics
- Expectation and Random Vectors
- Discrete Probability
- Example discrete distributions
- Continuous probability
- Gaussian Distribution
- Maximum Likelihood Estimation



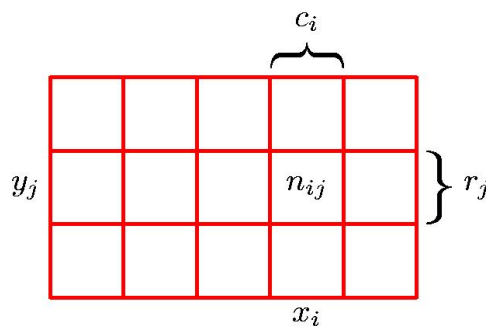
Think Generatively



Data generated from linear model

Want degree of **confidence** in **parameter values** and **predictions**

Probability Theory



•Marginal Probability

$$p(X = x_i) = \frac{c_i}{N}.$$

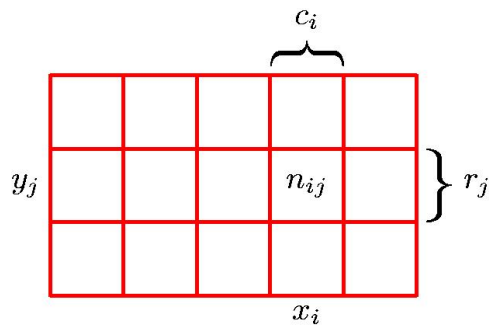
Joint Probability

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N}$$

•Conditional Probability

$$p(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i}$$

Probability Theory



•Sum Rule

$$p(X = x_i) = \frac{c_i}{N} = \frac{1}{N} \sum_{j=1}^L n_{ij}$$

$$= \sum_{j=1}^L p(X = x_i, Y = y_j)$$

Product Rule

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N} = \frac{n_{ij}}{c_i} \cdot \frac{c_i}{N}$$

$$= p(Y = y_j | X = x_i) p(X = x_i)$$

The Rules of Probability

- Sum Rule

$$p(X) = \sum_Y p(X, Y)$$

- Product Rule

$$p(X, Y) = p(Y|X)p(X)$$

Bayes' Theorem

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

$$p(X) = \sum_Y p(X|Y)p(Y)$$

posterior \propto likelihood \times prior

Expectation

The expected value of a function of a random variable X that is distributed according to $P(X)$

$$\mathbb{E}_{P(x)} \{f(X)\} = \sum_x f(x)P(x)$$

The expected value of the random variable X itself: the **mean**

$$\mathbb{E}_{P(x)} \{X\} = \sum_x xP(x)$$

Expectation

$$\mathbf{E}_{P(x)} \{f(X)\} = \sum_x f(x)P(x)$$

The expectation of the value of X if X is a fair die:

$$\mathbf{E}_{P(x)} \{X\} = \sum_x x \frac{1}{6} = \frac{1}{6} + \frac{2}{6} + \dots + \frac{6}{6} = \frac{21}{6} = 3.5$$

$$\mathbf{E}_{P(x)} \{X^2\} = \sum_x x^2 \frac{1}{6} = \frac{1}{6} + \frac{4}{6} + \dots + \frac{36}{6} = \frac{91}{6}$$

$$12.25 \neq 15.16$$

$$(\mathbf{E}_{P(x)} \{X\})^2 \neq \mathbf{E}_{P(x)} \{X^2\}$$

Expectation

$$\mathbf{E}_{P(x)} \{f(X)\} = \sum_x f(x)P(x)$$

$$(\mathbf{E}_{P(x)} \{X\})^2 \neq \mathbf{E}_{P(x)} \{X^2\}$$

In general: the expected value of a function of X is not equal to the function evaluated at the expected value of X !

usually

$$f(\mathbf{E}_{P(x)} \{X\}) \neq \mathbf{E}_{P(x)} \{f(X)\}$$

Special cases:

$$f(X) = aX$$

$$f(X) = a$$

$$\mathbf{E}_{P(x)} \{f(X) + g(X)\} = \mathbf{E}_{P(x)} \{f(X)\} + \mathbf{E}_{P(x)} \{g(X)\}$$

Expectation: Variance

$$\mathbf{E}_{P(x)} \{f(X)\} = \sum_x f(x)P(x)$$

Variance:

$$\text{var}\{X\} = \mathbf{E}_{P(x)} \{(X - \mathbf{E}_{P(x)} \{x\})^2\}$$

$$\begin{aligned} \text{var}\{X\} &= \mathbf{E}_{P(x)} \{(X - \mathbf{E}_{P(x)} \{x\})^2\} \\ &= \mathbf{E}_{P(x)} \{X^2 - 2X\mathbf{E}_{P(x)} \{X\} + \mathbf{E}_{P(x)} \{x\}^2\} \\ &= \mathbf{E}_{P(x)} \{X^2\} - 2\mathbf{E}_{P(x)} \{X\} \mathbf{E}_{P(x)} \{X\} + \mathbf{E}_{P(x)} \{X\}^2 \end{aligned}$$

$$\text{var}\{X\} = \mathbf{E}_{P(x)} \{X^2\} - \mathbf{E}_{P(x)} \{X\}^2$$

Vector Random Variables

Vector random variables!

$$p(\mathbf{x}) = p(x_1, x_2, \dots, x_N) = P(X_1 = x_1, X_2 = x_2, \dots, X_N = x_N)$$

Mean: $\mathbf{E}_{P(\mathbf{x})} \{\mathbf{x}\} = \sum_{\mathbf{x}} \mathbf{x}P(\mathbf{x})$

Very similar to *scalar* version:

$$\mathbf{E}_{P(x)} \{X\} = \sum_x xP(x)$$

Covariance:

$$\text{cov}\{\mathbf{x}\} = \mathbf{E}_{P(\mathbf{x})} \{(\mathbf{x} - \mathbf{E}_{P(\mathbf{x})} \{\mathbf{x}\})(\mathbf{x} - \mathbf{E}_{P(\mathbf{x})} \{\mathbf{x}\})^T\}$$

$$\begin{aligned} \text{cov}\{\mathbf{x}\} &= \mathbf{E}_{P(\mathbf{x})} \{(\mathbf{x} - \mathbf{E}_{P(\mathbf{x})} \{\mathbf{x}\})(\mathbf{x} - \mathbf{E}_{P(\mathbf{x})} \{\mathbf{x}\})^T\} \\ &= \mathbf{E}_{P(\mathbf{x})} \{\mathbf{x}\mathbf{x}^T - 2\mathbf{x}\mathbf{E}_{P(\mathbf{x})} \{\mathbf{x}\}^T + \mathbf{E}_{P(\mathbf{x})} \{\mathbf{x}\} \mathbf{E}_{P(\mathbf{x})} \{\mathbf{x}\}^T\} \end{aligned}$$

$$\text{cov}\{\mathbf{x}\} = \mathbf{E}_{P(\mathbf{x})} \{\mathbf{x}\mathbf{x}^T\} - \mathbf{E}_{P(\mathbf{x})} \{\mathbf{x}\} \mathbf{E}_{P(\mathbf{x})} \{\mathbf{x}\}^T$$

Vector Random Variables

Vector random variables!

$$p(\mathbf{x}) = p(x_1, x_2, \dots, x_N) = P(X_1 = x_1, X_2 = x_2, \dots, X_N = x_N)$$

Mean: $\mathbf{E}_{P(\mathbf{x})} \{\mathbf{x}\} = \sum_{\mathbf{x}} \mathbf{x} P(\mathbf{x})$

Very similar to **scalar** version:

$$\mathbf{E}_{P(x)} \{X\} = \sum_x x P(x)$$

When we move to vector random variables and consider their “variance”, the scalar version of variance needs to be extended...

Scalar variance:

$$\begin{aligned} \text{var}\{X\} &= \mathbf{E}_{P(x)} \{(X - \mathbf{E}_{P(x)} \{x\})^2\} \\ \text{var}\{X\} &= \mathbf{E}_{P(x)} \{X^2\} - \mathbf{E}_{P(x)} \{X\}^2 \end{aligned}$$

The scalar “summation” form of variance:

$$\begin{aligned} \text{var}(X) &= \sum_x (x - \mu_X)^2 \\ &= \sum_x (x - \mu_X)(x - \mu_X) \end{aligned}$$

This is like comparing a variable to itself

$$\text{var}\{X\} = \mathbf{E}_{P(x)} \{(X - \mathbf{E}_{P(x)} \{x\})(X - \mathbf{E}_{P(x)} \{x\})\}$$

When we want to calculate how one random variable (co)varies with another, Then we are interested in the **covariance**:

$$\text{cov}(X, Y) = \mathbf{E}_{p(x,y)} \{(x - \mathbf{E}_{p(x)} \{x\})(y - \mathbf{E}_{p(y)} \{y\})\}$$



13

(Co)variance of a Random Vector

- Covariance

$$\text{cov}(X, Y) = \mathbf{E}_{p(x,y)} \{(x - \mathbf{E}_{p(x)} \{x\})(y - \mathbf{E}_{p(y)} \{y\})\}$$

- Now, if we want to take the “variance” of a random vector, which is essentially a compact representation of a **joint distribution**, then we need to keep track of all of the pair-wise **covariances** of each of the random vector components, and we do this in the **covariance matrix**:

$$\Sigma = \begin{bmatrix} \mathbf{E}[(X_1 - \mu_1)(X_1 - \mu_1)] & \mathbf{E}[(X_1 - \mu_1)(X_2 - \mu_2)] & \cdots & \mathbf{E}[(X_1 - \mu_1)(X_n - \mu_n)] \\ \mathbf{E}[(X_2 - \mu_2)(X_1 - \mu_1)] & \mathbf{E}[(X_2 - \mu_2)(X_2 - \mu_2)] & \cdots & \mathbf{E}[(X_2 - \mu_2)(X_n - \mu_n)] \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{E}[(X_n - \mu_n)(X_1 - \mu_1)] & \mathbf{E}[(X_n - \mu_n)(X_2 - \mu_2)] & \cdots & \mathbf{E}[(X_n - \mu_n)(X_n - \mu_n)] \end{bmatrix}.$$

$$\text{cov}\{\mathbf{x}\} = \mathbf{E}_{P(\mathbf{x})} \{(\mathbf{x} - \mathbf{E}_{P(\mathbf{x})} \{\mathbf{x}\})(\mathbf{x} - \mathbf{E}_{P(\mathbf{x})} \{\mathbf{x}\})^T\}$$

$$\text{cov}\{\mathbf{x}\} = \mathbf{E}_{P(\mathbf{x})} \{\mathbf{x}\mathbf{x}^T\} - \mathbf{E}_{P(\mathbf{x})} \{\mathbf{x}\} \mathbf{E}_{P(\mathbf{x})} \{\mathbf{x}\}^T$$

cov{x} is shorthand for cov(x,x)

When x is a random vector, x,

then this is a matrix, Σ



14

Discrete Distributions

- The functions that characterize the discrete random variable are often referred to as *probability **mass** functions* (pmf)
- Bernoulli
- Binomial
- Multinomial

Bernoulli Distribution

- Coin flipping: heads=1, tails=0

$$p(x = 1|\mu) = \mu$$

- Bernoulli Distribution

$$\begin{aligned}\text{Bern}(x|\mu) &= \mu^x(1 - \mu)^{1-x} \\ \mathbb{E}[x] &= \mu \\ \text{var}[x] &= \mu(1 - \mu)\end{aligned}$$

Binomial Distribution

- N coin flips:

$$p(m \text{ heads} | N, \mu)$$

- Binomial Distribution

$$\text{Bin}(m|N, \mu) = \binom{N}{m} \mu^m (1 - \mu)^{N-m}$$

$$\mathbb{E}[m] \equiv \sum_{m=0}^N m \text{Bin}(m|N, \mu) = N\mu$$

$$\text{var}[m] \equiv \sum_{m=0}^N (m - \mathbb{E}[m])^2 \text{Bin}(m|N, \mu) = N\mu(1 - \mu)$$

Binomial Distribution

