

Detection

Binary Hypothesis Training:

- Two hypotheses H_0 and H_1
- If H_0 occurs, X has PMF $P_{X|H_0}(x)$
- If H_1 occurs, X has PMF $P_{X|H_1}(x)$
- $U_0 = \{\text{decide } H_0 \text{ based on } X\}$
- $U_1 = \{\text{decide } H_1 \text{ based on } X\}$
- Error occurs if we decide the wrong hypothesis based on X
- Goal is to minimize the **probability of error**:

$$P[\text{error}] = P[U_1|H_0]P[H_0] + P[U_0|H_1]P[H_1]$$

- Likelihood Ratio: $L(x) = \frac{P_{X|H_1}(x)}{P_{X|H_0}(x)}$.
- Log-Likelihood Ratio: $\ln(L(x)) = \ln\left(\frac{P_{X|H_1}(x)}{P_{X|H_0}(x)}\right)$.

Maximum Likelihood (ML) Rule:

- Intuition: Choose the hypothesis that best explains the observation.
- In terms of the conditional PMFs,
 - Decide H_0 if $P_{X|H_0}(x) > P_{X|H_1}(x)$.
 - Decide H_1 if $P_{X|H_0}(x) < P_{X|H_1}(x)$.
 - Decide either if $P_{X|H_0}(x) = P_{X|H_1}(x)$. *if diff*
- In terms of the likelihood ratio,
 - Decide H_0 if $L(x) < 1$.
 - Decide H_1 if $L(x) > 1$.
 - Decide either if $L(x) = 1$. *if the same*
- In terms of the log-likelihood ratio,
 - Decide H_0 if $\ln(L(x)) < 0$.
 - Decide H_1 if $\ln(L(x)) > 0$.
 - Decide either if $\ln(L(x)) = 0$. *with e*

Maximum a Posteriori (MAP) Rule:

- Intuition: Choose the most likely hypothesis given the observation.
- Attains the **minimum probability of error**.
- In terms of the conditional PMFs,
 - Decide H_0 if $P_{X|H_0}(x)P[H_0] > P_{X|H_1}(x)P[H_1]$.
 - Decide H_1 if $P_{X|H_0}(x)P[H_0] < P_{X|H_1}(x)P[H_1]$.
 - Decide either if $P_{X|H_0}(x)P[H_0] = P_{X|H_1}(x)P[H_1]$.
- In terms of the likelihood ratio,
 - Decide H_0 if $L(x) < \frac{P[H_0]}{P[H_1]}$.
 - Decide H_1 if $L(x) > \frac{P[H_0]}{P[H_1]}$.
 - Decide either if $L(x) = \frac{P[H_0]}{P[H_1]}$.
- In terms of the log-likelihood ratio,
 - Decide H_0 if $\ln(L(x)) < \ln\left(\frac{P[H_0]}{P[H_1]}\right)$.
 - Decide H_1 if $\ln(L(x)) > \ln\left(\frac{P[H_0]}{P[H_1]}\right)$.
 - Decide either if $\ln(L(x)) = \ln\left(\frac{P[H_0]}{P[H_1]}\right)$.

Estimation

- Two continuous random variables X and Y . Observe one and want to estimate other
- Here, we assume that we observe Y and estimate X using a function of Y , written as $x\hat{}$ (Y)
- Need to measure the error between the true X and our estimate $x\hat{}$ (Y).

- In this class, we focused on the **mean-squared error (MSE)**:

$$\text{MSE} = E[(X - \hat{x}(Y))^2]$$

- **ML Estimator:** $\hat{x}_{\text{ML}}(y) = \arg \max_x f_{Y|X}(y|x)$
- **MAP Estimator:** $\hat{x}_{\text{MAP}}(y) = \arg \max_x f_{Y|X}(y|x)f_X(x)$

Minimum Mean-Squared Error (MMSE) Estimator:

- Attains the minimum MSE amongst all possible estimators
- Given by the conditional expectation:

$$\hat{x}_{\text{MMSE}}(y) = E[X|Y = y]$$
- **Unbiased:** $E[\hat{x}_{\text{MMSE}}(Y)] = E[X]$
- Orthogonality Principle: $E[(X - \hat{x}_{\text{MMSE}}(Y))g(y)] = 0$ for any function $g(y)$ (including $\hat{x}_{\text{MMSE}}(Y)$).

Linear Least-Squares Error (LLSE) Estimator:

- Attains the **minimum MSE** amongst all linear estimators.
- Two formulas: $\hat{x}_{\text{LLSE}}(y) = E[X] + \frac{\text{Cov}[X, Y]}{\text{Var}[Y]}(y - E[Y])$

$$\hat{x}_{\text{LLSE}}(y) = E[X] + \rho_{X,Y} \frac{\sigma_X}{\sigma_Y}(y - E[Y])$$
- Unbiased: $E[\hat{x}_{\text{LLSE}}(Y)] = E[X]$.
- Orthogonality Principle: $E[(X - \hat{x}_{\text{LLSE}}(Y))(aY + b)] = 0$ for any linear function $ay + b$ (including $\hat{x}_{\text{LLSE}}(Y)$).

Sums of Random Variables

More than Two Random Variables:

- Joint PMF for n discrete random variables $P_{X_1, \dots, X_n}(x_1, \dots, x_n)$
- Joint PDF for n continuous random variables $f_{X_1, \dots, X_n}(x_1, \dots, x_n)$
- Expected Value of a Sum: $E[X_1 + \dots + X_n] = \sum_{i=1}^n E[X_i]$
- Variance of a Sum: $\text{Var}[X_1 + \dots + X_n] = \sum_{i=1}^n \sum_{j=1}^n \text{Cov}[X_i, X_j]$

Independent and Identically Distributed Random Variables:

- Joint PMF $P_{X_1, \dots, X_n}(x_1, \dots, x_n) = P_X(x_1) \cdots P_X(x_n)$
- Joint PDF $f_{X_1, \dots, X_n}(x_1, \dots, x_n) = f_X(x_1) \cdots f_X(x_n)$
- Expected Value of a Sum: $E[X_1 + \dots + X_n] = nE[X]$
- Variance of a Sum: $\text{Var}[X_1 + \dots + X_n] = n\text{Var}[X]$

Limit Theorems

- **Markov's Inequality:** For any non-negative random variable X and constant $c > 0$,

$$P[X \geq c] \leq \frac{E[X]}{c}$$
- **Chebyshev's Inequality:** For any random variable X and constant $c > 0$,

$$P[|X - E[X]| \geq c] \leq \frac{\text{Var}[X]}{c^2}$$
- **Weak Law of Large Numbers:** Let X_1, X_2, \dots, X_n be i.i.d. random variables. For any constant $c > 0$,

$$P\left(\left|\frac{1}{n} \sum_{i=1}^n X_i - E[X]\right| \geq c\right) \leq \frac{\text{Var}[X]}{nc^2}$$

- **Hoeffding's Inequality:** Let X_1, X_2, \dots, X_n be i.i.d. Bernoulli(p) random variables. For any constant $c > 0$,

$$P\left(\left|\frac{1}{n} \sum_{i=1}^n X_i - E[X]\right| \geq c\right) \leq 2e^{-2nc^2}$$

Intro to Machine Learning

- The goal is to decide between two hypotheses, but we do not have access to the underlying probability model
- Instead we have a dataset D consisting of n sample $\mathcal{D} = \{(\vec{X}_1, Y_1), (\vec{X}_2, Y_2), \dots, (\vec{X}_n, Y_n)\}$

The i^{th} sample (X_i, Y_i) has an observation vector X_i and a label Y_i which we assume is -1 or +1

- We use this dataset to come up with a classifier $h(X)$, which is a function that maps any possible observation vector into a guess of its label
- Ideally, given enough samples, we would approach the probability of error of the optimal MAP rule that knows the probability distribution

Training and test Error:

- To make sure we are not overfitting, we split our dataset into non-overlapping training and test datasets:

$$\mathcal{D}_{\text{train}} = \{(\vec{X}_{\text{train},1}, Y_{\text{train},1}), (\vec{X}_{\text{train},2}, Y_{\text{train},2}), \dots, (\vec{X}_{\text{train},n_{\text{train}}}, Y_{\text{train},n_{\text{train}}})\}$$

$$\mathcal{D}_{\text{test}} = \{(\vec{X}_{\text{test},1}, Y_{\text{test},1}), (\vec{X}_{\text{test},2}, Y_{\text{test},2}), \dots, (\vec{X}_{\text{test},n_{\text{test}}}, Y_{\text{test},n_{\text{test}}})\}.$$

- The training set is used to construct our classifier $h(X)$ and the test set can only be used to evaluate its performance.

$$\text{Training Error} = \frac{1}{n_{\text{train}}} \sum_{i=1}^{n_{\text{train}}} g_{\text{not_equal}}(h(\vec{X}_{\text{train},i}), Y_{\text{train},i}),$$

$$\text{Test Error} = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} g_{\text{not_equal}}(h(\vec{X}_{\text{test},i}), Y_{\text{test},i})$$

$$\text{where } g_{\text{not_equal}}(y_{\text{guess}}, y_{\text{true}}) = \begin{cases} 1 & \text{if } y_{\text{guess}} \neq y_{\text{true}} \\ 0 & \text{if } y_{\text{guess}} = y_{\text{true}} \end{cases}$$

Basic Classifiers:

- The **closest average** classifier first computes the average vector for each label, based on the training set. Given a new observation vector, it computes the distance to each average and choose the label with the smallest distance
- The **nearest neighbor** classifier, when given a new observation vector, computes the distance to every sample in the training set to find the closest point. It then outputs the label of this point as its guess
- The **linear regression** classifier first computes a weighted sum of a new observation vector, and outputs the sign of this weighted sum and its guess. The training set is used to select the weights

Dimensionality Reduction :

- Principal component analysis** allows us to reduce the dimensionality of our observations, by only keeping the (orthogonal) directions corresponding to the largest variance

Markov Chain

- Sequence of (discrete) random variables X_0, X_1, X_2, \dots such that, given the history X_0, \dots, X_n , the next state X_{n+1} only depends on the current state X_n . This is sometimes called the **Markov Property**:

$$P[X_{n+1} = j | X_n = i, X_{n-1} = i_{n-1}, \dots, X_0 = i_0] = P_{ij} \quad \text{for all } n$$

- The **transition probabilities** P_{ij} are the probabilities of moving from state i to state j in one time step
- The n -step transition probabilities $P_{ij}(N)$ are the probabilities of moving from state i to state j in exactly n time steps. They can be determined via the **Chapman-Kolmogorov equations**:

$$P_{ij}(n+m) = \sum_{k=0}^K P_{ik}(n) P_{kj}(m)$$

State Transition Matrix:

- The state transition matrix is:

$$; \mathbf{P} = \begin{bmatrix} P_{00} & P_{01} & \cdots & P_{0K} \\ P_{10} & P_{11} & \cdots & P_{1K} \\ \vdots & \vdots & \ddots & \vdots \\ P_{K0} & P_{K1} & \cdots & P_{KK} \end{bmatrix}$$

- Row index is for the current state, column index is for the next state
- All rows must sum to 1

Probability State Vector:

- The probability state vector $\vec{p}(n) = [p_0(n) \ p_1(n) \ \cdots \ p_K(n)]$ captures the PMF across the states at time n .
- $p_i(n) = P[X_n = i]$
- The entries of $\vec{p}(n)$ must sum to 1.
- Moving forward one time step: $\vec{p}(n+1) = \vec{p}(n)\mathbf{P}$.
- Moving forward n time steps: $\vec{p}(n) = \vec{p}(0)\mathbf{P}^n$.

State Classification:

- State j is **accessible** from state i if it is possible to reach state j starting from state i in one or more time steps. Notation: $i \rightarrow j$
- States i and j **communicate** if $i \rightarrow j$ and $j \rightarrow i$. Notation: $i \leftrightarrow j$
- A **communicating class C** is a subset of states such that if $i \in C$, then $j \in C$ if and only if $i \leftrightarrow j$.
- A Markov chain is **irreducible** if all of its states belong to a single communicating class.
- A state j is **transient** if there is a state k such that $j \rightarrow k$ but $k \not\rightarrow j$.
- Any state that is not transient is **recurrent**.
- The **period d** of a state i is the greatest common divisor of the length of all cycles from i back to itself.
- If the period is 1, then the state is called **aperiodic**. The entire Markov chain is aperiodic if all states are aperiodic.

Limiting State Probability Vector:

- For an irreducible, aperiodic Markov chain, there is a unique limiting state probability vector

$$\vec{\pi} = \lim_{n \rightarrow \infty} \vec{p}(n)$$

Properties of the Limiting State Probability Vector:

- Normalization: $\sum_{i=0}^K \pi_i = 1$
- Any initial state probability vector $\vec{p}(0)$ will converge to $\vec{\pi}$.
- Steady-State Distribution: $\vec{\pi} = \vec{\pi} \mathbf{P}$.

Handling Transient States:

- If there is only one recurrent communicating class, then there is still a unique limiting state probability vector

$$\vec{\pi} = \lim_{n \rightarrow \infty} \vec{p}(n).$$

- If state j is transient, then $\pi_j = 0$

LAST EXAM STUFF

Important Families of Discrete Random Variables:

Bernoulli Random Variables:

- X is a Bernoulli(p) random variable if it has PMF

$$P_X(x) = \begin{cases} 1-p & x = 0, \\ p & x = 1. \end{cases}$$

- Range: $S_X = \{0, 1\}$.
- Expected Value: $E[X] = p$.
- Variance: $\text{Var}[X] = p(1-p)$.
- Interpretation: Single trial with success probability p .

Geometric Random Variables:

- X is a Geometric(p) random variable if it has PMF

$$P_X(x) = \begin{cases} p(1-p)^{x-1} & x = 1, 2, \dots, \\ 0 & \text{otherwise.} \end{cases}$$

- Range: $S_X = \{1, 2, \dots\}$.
- Expected Value: $E[X] = \frac{1}{p}$.
- Variance: $\text{Var}[X] = \frac{1-p}{p^2}$.
- Interpretation: # of independent Bernoulli(p) trials until first success.

Binomial Random Variables:

- X is a Binomial(n, p) random variable if it has PMF

$$P_X(x) = \begin{cases} \binom{n}{x} p^x (1-p)^{n-x} & x = 0, 1, \dots, n, \\ 0 & \text{otherwise.} \end{cases}$$

- Range: $S_X = \{0, 1, \dots, n\}$.
- Expected Value: $E[X] = np$.
- Variance: $\text{Var}[X] = np(1-p)$.
- Interpretation: # of successes in n independent Bernoulli(p) trials.

Discrete Uniform Random Variables:

- X is a Discrete Uniform(k, ℓ) random variable if it has PMF

$$P_X(x) = \begin{cases} \frac{1}{\ell - k + 1} & x = k, k+1, \dots, \ell, \\ 0 & \text{otherwise.} \end{cases}$$

- Range: $S_X = \{k, k+1, \dots, \ell\}$.
- Expected Value: $E[X] = \frac{k+\ell}{2}$.
- Variance: $\text{Var}[X] = \frac{(\ell-k)(\ell-k+2)}{12} = \frac{(\ell-k+1)^2 - 1}{12}$.
- Interpretation: equally likely to take any integer value from k to ℓ .

Poisson Random Variables:

- X is a Poisson(α) random variable if it has PMF

$$P_X(x) = \begin{cases} \frac{\alpha^x}{x!} e^{-\alpha} & x = 0, 1, \dots, \\ 0 & \text{otherwise.} \end{cases}$$

- Range: $S_X = \{0, 1, \dots\}$.
- Expected Value: $E[X] = \alpha$.
- Variance: $\text{Var}[X] = \alpha$.
- Interpretation: # of arrivals in a fixed time window.

Important Families of Continuous Random Variables:

Uniform Random Variables:

- X is a Uniform(a, b) random variable if it has PDF

$$f_X(x) = \begin{cases} \frac{1}{b-a} & a \leq x < b \\ 0 & \text{otherwise.} \end{cases}$$

- CDF:

$$F_X(x) = \begin{cases} 0 & x < a \\ \frac{x-a}{b-a} & a \leq x < b \\ 1 & b \leq x \end{cases}$$

- Expected Value: $E[X] = (a+b)/2$
- Variance: $\text{Var}[X] = (b-a)^2/12$

Exponential Random Variables:

- X is an Exponential(λ) random variable if it has PDF

$$f_X(x) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0 \\ 0 & x < 0 \end{cases}$$

- CDF:

$$F_X(x) = \begin{cases} 1 - e^{-\lambda x} & x \geq 0 \\ 0 & x < 0 \end{cases}$$

- Expected Value: $E[X] = 1/\lambda$
- Variance: $\text{Var}[X] = 1/\lambda^2$

Gaussian Random Variables:

- X is an Gaussian(μ, σ^2) random variable if it has PDF

$$f_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

- CDF:

$$F_X(x) = \Phi\left(\frac{x-\mu}{\sigma}\right) \text{ where } \Phi(z) = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-\frac{w^2}{2}} dw$$

- $\Phi(z)$ is the standard normal CDF and $Q(z) = 1 - \Phi(z)$ is the standard normal complementary CDF
- $\Phi(-z) = 1 - \Phi(z) = Q(z)$
- Expected Value: $E[X] = \mu$
- Variance: $\text{Var}[X] = \sigma^2$

- $\text{Var}[X] = E[X^2] - E[X]^2$
- $\text{Var}[aX+b] = a^2 \text{Var}[X]$
- $\sigma_X = \sqrt{\text{Var}[X]}$
- $E[X] = \int x f_X(x) dx$

- The covariance of random variables X and Y is

$$\text{Cov}[X, Y] = E[(X - E[X])(Y - E[Y])]$$

- Captures how X and Y vary together (linearly).
- Another useful formula is $\text{Cov}[X, Y] = E[XY] - E[X]E[Y]$
- Sometimes, $E[XY]$ is called the correlation. Don't confuse this with the correlation coefficient (next slide).

Covariance of Linear Functions: For any constants a, b, c , and d ,

$$\text{Cov}[aX + b, cY + d] = ac \text{Cov}[X, Y].$$

Variance of Sums:

$$\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y] + 2\text{Cov}[X, Y].$$

Variance of Weighted Sums:

$$\text{Var}[aX + bY + c] = a^2 \text{Var}[X] + b^2 \text{Var}[Y] + 2ab\text{Cov}[X, Y].$$

- The correlation coefficient is $\rho_{X,Y} = \frac{\text{Cov}[X, Y]}{\sqrt{\text{Var}[X]\text{Var}[Y]}}$

$$\bullet \quad P[A \cap B] = P[A]P[B|A] = P[B]P[A|B].$$

- The events A and B are conditionally independent given C if

$$P[A \cap B|C] = P[A|C]P[B|C].$$

Pairs of Continuous Random Variables:

- Joint PDF: $f_{X,Y}(x, y) = \frac{\partial^2 F_{X,Y}(x, y)}{\partial x \partial y}$
- Range $S_{X,Y} = \{(x, y) : f_{X,Y}(x, y) > 0\}$.
- Marginal PDFs $f_X(x)$ and $f_Y(y)$ are just the PMFs of the individual random variables X and Y , respectively.

$$f_X(x) = \int_{-\infty}^{\infty} f_{X,Y}(x, y) dy \quad f_Y(y) = \int_{-\infty}^{\infty} f_{X,Y}(x, y) dx$$

- Conditional PDFs give the probability density of one random variable when the other is fixed to a certain value:

$$f_{X|Y}(x|y) = \begin{cases} \frac{f_{X,Y}(x,y)}{f_Y(y)} & f_Y(y) > 0 \\ 0 & \text{otherwise.} \end{cases}$$

$$f_{Y|X}(y|x) = \begin{cases} \frac{f_{X,Y}(x,y)}{f_X(x)} & f_X(x) > 0 \\ 0 & \text{otherwise.} \end{cases}$$

EXAMPLE PROBLEMS

T/F

- (a) Let A and B be events. Then, $P[B|A] = P[A|B] \frac{1 - P[B^c]}{1 - P[A^c]}$.

True. By Bayes' Theorem and the complement property,

$$P[B|A] = \frac{P[A|B]P[B]}{P[A]} = P[A|B] \frac{1 - P[B^c]}{1 - P[A^c]}$$

- (b) Let X and Y be a pair of uncorrelated random variables and let $Z = 2X - Y$. Then, $\text{Var}[Z] = 4\text{Var}[X] - \text{Var}[Y]$.

False. $\text{Var}[2X - Y] = 4\text{Var}[X] + \text{Var}[Y]$.

- (c) Let X and Y be jointly Gaussian random variables with $E[X] = E[Y] = 1$, $\text{Var}[X] = \text{Var}[Y] = 5$, and $\rho_{X,Y} = 0$. Also, let $Z = X + Y$. Then, $P[Z \geq 2] = 1/2$.

True. Since Z is a linear function of jointly Gaussian random variables, it is also Gaussian. It has mean $E[Z] = E[X + Y] = E[X] + E[Y] = 2$. Since the Gaussian PDF is symmetric about its mean, $P[Z \geq 2] = 1/2$.

- (d) Suppose we have a Markov chain with two communicating classes. Then, there is no unique limiting state probability vector.

False. If one of the classes is transient, there will still be a unique limiting state probability vector.

- (e) Let X be a zero-mean random variable with variance 1. Then, $P[X^2 \geq 2] \leq 1/2$.

True. Define $Y = X^2$ and note that $E[Y] = E[X^2] = \text{Var}[X] + (E[X])^2 = 1$. Using Markov's inequality, $P[Y \geq 2] \leq \frac{E[Y]}{2} = 1/2$.

- (f) Let X be a non-negative random variable and $Y = \ln(X)$. Then, $F_Y(y) = F_X(e^y)$.

True. $F_Y(y) = P[Y \leq y] = P[\ln(X) \leq y] = P[X \leq e^y] = F_X(e^y)$.

- (a) Let X and Y be independent random variables, each with mean 1 and variance 1. Then, $E[X^2 Y^2] = 4$.

True. $E[X^2] = \text{Var}[X] + (E[X])^2 = 1 + 1 = 2$, and $E[Y^2] = \text{Var}[Y] + (E[Y])^2 = 1 + 1 = 2$. Using independence, $E[X^2 Y^2] = E[X^2]E[Y^2]$, so the result follows.

- (b) Let X and Y be continuous random variables. Then, for any constants $a > 0$ and $b > 0$, $P[\{X > a\} \cap \{Y > b\}] = \left(\int_a^{\infty} f_X(x) dx \right) \cdot \left(\int_b^{\infty} f_Y(y) dy \right)$

False. The probability of interest can be written as $\int_a^{\infty} \int_b^{\infty} f_{X,Y}(x, y) dy dx$, but this cannot be simplified to the given expression unless X and Y are independent.

- (c) Let X and Y be jointly Gaussian random variables with $E[X] = E[Y] = 0$, $\text{Var}[X] = \text{Var}[Y] = 1$, and $\rho_{X,Y} = 0$. Also, let $U = X + 2Y$ and $V = 2Y - X$. Then, U and V are independent.

False. Note that $E[U] = 0$, $E[V] = 0$, and

$$E[UV] = E[(X + 2Y)(2Y - X)] = E[4Y^2 - X^2] = 4 - 1 = 3,$$

so $\text{Cov}[U, V] = E[UV] - E[U]E[V] = 3 \neq 0$. Correlated Gaussian random variables are not independent.

- (d) Let X and Y be a pair of uncorrelated random variables ($\text{Cov}[X, Y] = 0$). Then, $E[(X - \hat{x}_{\text{LLSE}}(Y))^2] = \text{Var}[X]$.

True. Since $\text{Cov}[X, Y] = 0$, $\hat{x}_{\text{LLSE}}(Y) = E[X]$. Substituting this, we are left with the definition of variance.

- (e) For a Markov chain and any states i and j , $P[X_2 = j | X_0 = i] = P[X_2 = j | X_1 = i]$.

False. A two-step transition probability is not generally the same as a one-step transition probability.

- (f) For the MAP decision rule, $P[\text{error}] \leq \min(P[H_0], P[H_1])$.

True. $P[\text{error}] = P[H_0]$ is achieved by always deciding H_1 , and $P[\text{error}] = P[H_1]$ is achieved by always deciding H_0 , so the MAP decision rule (which minimizes the probability of error) will do at least as well as the minimum.

- (g) If X and Y are independent, $\hat{x}_{\text{MMSE}}(y) = 0$.

False. X and Y being independent makes $\hat{x}_{\text{MMSE}}(y) = E[X]$, which is not necessarily zero.

- (h) If $E[X] = E[Y] = 0$, then $E[XY] \leq \sqrt{\text{Var}[X]\text{Var}[Y]}$.

True. This follows from $\rho_{X,Y} \leq 1$.

6.1 Binary Hypothesis Testing

$$P_e = P[\{\text{error}\}] = P[\{\text{error}\} | H_0] P[H_0] + P[\{\text{error}\} | H_1] P[H_1]$$

$$= P[\{y \in A, \bar{y} | H_0\}] P[H_0] + P[\{y \in A_0, \bar{y} | H_1\}] P[H_1]$$

Error =

$$\begin{array}{ll} \text{Discrete} & \text{Continuous} \\ P[\{y \in A, \bar{y} | H_0\}] = \sum_{y \in A_1} P_{y|H_0}(y) & P[\{y \in A, \bar{y} | H_0\}] = \int_{A_1} f_{y|H_0}(y) dy \end{array}$$

Probability of false alarm: $P_{FA} = P[\{y \in A, \bar{y} | H_0\}]$

Probability of missed detection: $P_{MD} = P[\{y \in A_0, \bar{y} | H_1\}]$

$$\text{so: } P_e = P_{FA} P[H_0] + P_{MD} P[H_1]$$

Maximum Likelihood (ML) Rule: selects H w/ highest likelihood value

Discrete:

$$D^{\text{ML}} = \begin{cases} 1, & P_{y|H_1}(y) \geq P_{y|H_0}(y) \\ 0, & P_{y|H_1}(y) < P_{y|H_0}(y) \end{cases}$$

Continuous:

$$D^{\text{ML}} = \begin{cases} 1, & f_{y|H_1}(y) \geq f_{y|H_0}(y) \\ 0, & f_{y|H_1}(y) < f_{y|H_0}(y) \end{cases}$$

Maximum a posteriori (MAP) rule: Selects H that is most likely given observation

Discrete:

$$D^{\text{MAP}} = \begin{cases} 1, & P_{y|H_1}(y) P[H_1] \geq P_{y|H_0}(y) P[H_0] \\ 0, & P_{y|H_1}(y) P[H_1] < P_{y|H_0}(y) P[H_0] \end{cases}$$

Continuous:

$$D^{\text{MAP}}(y) = \begin{cases} 1, & f_{y|H_1}(y) P[H_1] \geq f_{y|H_0}(y) P[H_0] \\ 0, & f_{y|H_1}(y) P[H_1] < f_{y|H_0}(y) P[H_0] \end{cases}$$

* MAP Rule is optimal

6.2: Likelihood Ratio

Likelihood Ratio $L(y)$:

$$L(y) = \begin{cases} \frac{P_{y|H_1}(y)}{P_{y|H_0}(y)} & y \text{ is discrete} \\ \frac{f_{y|H_1}(y)}{f_{y|H_0}(y)} & y \text{ is continuous} \end{cases}$$

$$D^{\text{ML}}: \begin{cases} 1, & L(y) \geq 1 \\ 0, & L(y) < 1 \end{cases} = \begin{cases} 1, & \ln(L(y)) \geq 0 \\ 0, & \ln(L(y)) < 0 \end{cases}$$

$$D^{\text{MAP}}: \begin{cases} 1, & L(y) \geq \frac{P[H_1]}{P[H_0]} \\ 0, & L(y) < \frac{P[H_0]}{P[H_1]} \end{cases} = \begin{cases} 1, & \ln(L(y)) \geq \ln\left(\frac{P[H_1]}{P[H_0]}\right) \\ 0, & \ln(L(y)) < \ln\left(\frac{P[H_0]}{P[H_1]}\right) \end{cases}$$

* Remember for Gaussian:

$$F_x(x) = \Phi\left(\frac{x-\mu}{\sigma}\right) \quad * \text{Standard normal CDF}$$

$$Q(z) = 1 - \Phi(z) = \Phi(-z)$$

Probability of detection

$$P_D = P[\{y \in A, \bar{y} | H_1\}] = 1 - P[\{y \in A_0, \bar{y} | H_1\}] = 1 - P_{H_1}$$

ROC Curve:

① Find $L(y)$

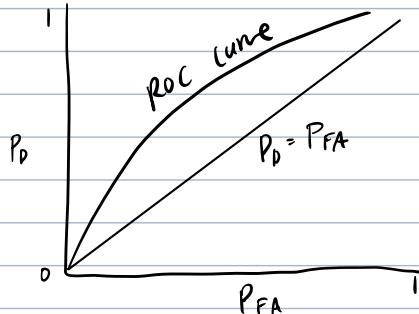
② Create a likelihood ratio test w/ threshold

$$T: D_T(y) = \begin{cases} 1, & L(y) \geq T \\ 0, & L(y) < T \end{cases}$$

③ Starting from $T=0$ to $T=\infty$, calculate

$$A_1(T), P_{FA}(T), P_D(T)$$

④ Plot resulting $(P_{FA}(T), P_D(T))$ pairs



properties:

- Concave

- Lies on or above

$$P_D = P_{FA} \text{ line}$$

6.3 Detection w/ unequal costs

Minimum Bayes Risk Detection:

\rightarrow cost C_{ij}

- C_{10} = cost of false alarm

- C_{01} = cost of missed detection

Expected Risk:

$$\text{Discrete: } E[R] = \sum_{y \in A_0} C_{01} P_{y|H_1}(y) P[H_1] + \sum_{y \in A_1} C_{10} P_{y|H_0}(y) P[H_0]$$

$$\text{Continuous: } E[R] = \int_{A_0} C_{01} f_{y|H_1}(y) P[H_1] dy + \int_{A_1} C_{10} f_{y|H_0}(y) P[H_0] dy$$

Minimum Baye's Risk decision Rule:

$$D^{MBR}(y) = \begin{cases} 1, & L(y) \geq \frac{C_{10}}{C_{01}} P[H_1] \\ 0, & L(y) < \frac{C_{10}}{C_{01}} P[H_1] \end{cases}$$

M-ary Hypotheses: H_0, H_1, \dots, H_{M-1}

$$P_e = \sum_{i=0}^{M-1} P[\{y \in A_i\} | H_i] P[H_i]$$

$$D^{ML} = \arg \max_{i \in \{0, 1, \dots, M-1\}} P_{y|H_i}(y) \quad (\text{discrete})$$

$$D^{ML} = \arg \max_{i \in \{0, 1, \dots, M-1\}} f_{y|H_i}(y) \quad (\text{continuous})$$

$$\star \arg \max (-g(y)) = \arg \min(g(y))$$

$$D^{MAP} = \arg \max_{i \in \{0, 1, \dots, M-1\}} P_{y|H_i}(y) P[H_i] \leftarrow$$

$$D^{MAP} = \arg \max_{i \in \{0, 1, \dots, M-1\}} f_{y|H_i}(y) P[H_i] \leftarrow \text{Requires knowing } P[H_i]$$

7.1: ML, MAP, and MMSE Estimation

Mean-Squared Error: X is unobserved, Y is observed

$$MSE = E[(X - \hat{x}(y))^2] \leftarrow \text{quality of prediction}$$

$\hat{x}(y)$ = estimation rule

Estimator $\hat{x}(y)$ is unbiased if the error $X - \hat{x}(y)$ has zero mean: $E[X - \hat{x}(y)] = 0$

ML Estimator: $\hat{x}_{ML}(y)$

$$Y \text{ discrete: } \hat{x}_{ML} = \arg \max_{x \in R_x} P_{y|x}(y|x)$$

$$Y \text{ continuous: } \hat{x}_{ML}(y) = \arg \max_{x \in R_x} f_{y|x}(y|x)$$

MAP estimator: $\hat{x}_{MAP}(y)$

$$X, Y \text{ discrete: } \arg \max_{x \in R_x} P_{y|x}(y|x) P_x(x) = \arg \max_{x \in R_x} P_{x|y}(x|y)$$

$$X, Y \text{ Jointly continuous: } \arg \max_{x \in R_x} f_{y|x}(y|x) f_x(x) = \arg \max_{x \in R_x} f_{x|y}(x|y)$$

$$X \text{ Discrete, } Y \text{ Continuous: } \arg \max_{x \in R_x} f_{y|x}(y|x) P_x(x)$$

$$X \text{ Continuous, } Y \text{ Discrete: } \arg \max_{x \in R_x} P_{y|x}(y|x) f_x(x)$$

$$\text{Unbiased? } E[\hat{x}_{ML}(y)] = E[E[\hat{x}_{ML}(y) | X]] = E[X] \quad \checkmark \text{ yes}$$

* In general, ML & MAP estimates are biased

MMSE estimator: OPTIMAL

* Attains smallest possible MSE

$$\hat{x}_{MMSE}(y) = E[x | y=y] \rightarrow \int_{-\infty}^{\infty} x f_{x|y}(x|y) dx$$

- unbiased

- Error of MMSE estimator is orthogonal to any function $g(y)$ of the observation

$$E[(X - \hat{x}_{MMSE}(y)) g(y)] = 0$$

Jointly Gaussian Random Vars

If X and Y are jointly Gaussian, MMSE estimator is a linear function of Y

$$\hat{X}_{\text{MMSE}}(y) = \mathbb{E}[X|Y=y] = \mu_x + \rho_{x,y} \frac{\sigma_x}{\sigma_y} (y - \mu_y) = \mathbb{E}[X] + \frac{\text{Cov}[X,Y]}{\text{Var}[Y]} (y - \mathbb{E}[Y])$$

$$\text{MSE}_{\text{MMSE}} = \mathbb{E}[(X - \hat{X}_{\text{MMSE}}(y))^2] = (1 - \rho_{x,y}^2) \sigma_x^2 = \text{Var}[X] - \frac{(\text{Cov}[X,Y])^2}{\text{Var}[Y]}$$

7.2: LLSE and Vector Estimation

The linear least squares error $\hat{X}_{\text{LLSE}}(y)$ attains the smallest possible mean-squared error among all linear estimators

$$\hat{X}_{\text{LLSE}}(y) = \mu_x + \rho_{x,y} \frac{\sigma_x}{\sigma_y} (y - \mu_y) = \mathbb{E}[X] + \frac{\text{Cov}[X,Y]}{\text{Var}[Y]} (y - \mathbb{E}[Y])$$

$$\text{MSE}_{\text{LLSE}} = \sigma_x^2 (1 - \rho_{x,y}^2) = \text{Var}[X] - \frac{(\text{Cov}[X,Y])^2}{\text{Var}[Y]}$$

* Unbiased: $\mathbb{E}[\hat{X}_{\text{LLSE}}(y)] = \mathbb{E}[X]$

* Error is orthogonal to any linear function $ay+b$: $\mathbb{E}[(X - \hat{X}_{\text{LLSE}}(y))(ay+b)] = 0$

$$\hat{X}_{\text{LLSE}}(y) = \mathbb{E}[X] + \sum_{x,y} \sum_{i=1}^{n-1} (y - \mathbb{E}[Y])$$

Cross Covariance matrix X, Y : $\Sigma_{x,y} = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])^T] = \begin{bmatrix} \text{Cov}[X_1, Y_1] & \dots & \text{Cov}[X_1, Y_n] \\ \vdots & \ddots & \vdots \\ \text{Cov}[X_n, Y_1] & \dots & \text{Cov}[X_n, Y_n] \end{bmatrix}$

Recall LLSE estimator Properties:

* unbiased

$$* \mathbb{E}[(X - \hat{X}_{\text{LLSE}}(Y))(AY + b)^T] = 0 \leftarrow \text{all zeros matrix}$$

8.1 Sums of Random Variables: Bounds

Random variables X, X_2, \dots, X_n are independent and identically distributed (i.i.d.) if they are independent and have the same PMF and PDF

Discrete PMF: $\prod_{i=1}^n P_{x_i}(x_i)$ Continuous PDF: $\prod_{i=1}^n f_{x_i}(x_i)$

$$\cdot \mathbb{E}[S_n] = \mathbb{E}\left[\sum_{i=1}^n X_i\right] = n\mathbb{E}[X]$$

$$\cdot \text{Var}[S_n] = \text{Var}\left[\sum_{i=1}^n X_i\right] = n\text{Var}[X] \quad * \text{compute using marginal PMF } P_x(x_i) \text{ or PDF } f_x(x_i)$$

$$M_n = \frac{1}{n} \sum_{i=1}^n X_i$$

$$\mathbb{E}[M_n] = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \mathbb{E}[X] \quad \text{Var}[M_n] = \text{Var}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n} \text{Var}[X]$$

Markov's Inequality:

$$\mathbb{P}[S_n \geq c] \leq \mathbb{E}[S_n]/c$$

Chebyshov's Inequality:

$$\mathbb{P}[|X - \mathbb{E}[X]| \geq c] \leq \text{Var}[X]/c^2$$

Hoeffding's Inequality

$$\mathbb{P}\left[|\sum_{i=1}^n X_i - \mathbb{E}\left[\sum_{i=1}^n X_i\right]| \geq c\right] \leq 2 \exp\left(\frac{-2c^2}{n(b-a)^2}\right)$$

8.2: Sums of Random Variables: Limits

Sample Mean: $M_n = \frac{1}{n} \sum_{i=1}^n X_i$

Weak Law of Large Numbers: $\lim_{n \rightarrow \infty} \mathbb{P}[|M_n - \mu| > \varepsilon] = 0$

Strong Law of Large Numbers: $\mathbb{P}\left[\lim_{n \rightarrow \infty} M_n = \mu\right] = 1$

Central Limit Theorem: $\lim_{n \rightarrow \infty} F_{Y_n}(y) = \Phi(y)$ where $Y_n = \frac{S_n - n\mu}{\sigma\sqrt{n}}$ * sum of many small independent effect look gaussian eventually

CDF for Gaussian(0,1)

9.1 Confidence Intervals

Sample Variance $V_n = \frac{1}{n-1} \sum_{i=1}^n (X_i - M_n)^2$ *unbiased estimator of variance

· Confidence level: $1-\alpha$ satisfies

$$P[A \leq \mu \leq B] = 1-\alpha \quad (\text{confidence interval: } [A, B])$$

Assumptions:

① Known Variance:

Confidence Interval: $[M_n - \varepsilon, M_n + \varepsilon]$

confidence level: $1 - \theta^2/n\varepsilon^2$

② CLT and known Variance:

$$CI: [M_n + \frac{\sigma}{\sqrt{n}}, M_n + \frac{\sigma}{\sqrt{n}}]$$

$$CL = 1-\alpha$$

*find γ SD so that $Q(\gamma) = \alpha/2$

$$\star \varepsilon = \sigma/\sqrt{n}$$

* always round up

Chi-square RV w/ n degrees-of-freedom =

$Y \sim \chi^2(n)$: $Y = \sum_{i=1}^n X_i^2$ where $X_i \sim N(0,1)$
(sum of squares of n independent standard Gaussian RVs)

→ Mean: $E[Y] = n$

→ Var: $\text{Var}[Y] = 2n$ CDF: look up table

Student's t-distribution with n-degrees of freedom:

$W \sim T(n)$: where $Z \sim N(0,1)$, $Y \sim \chi^2(n)$ and Y, Z independent

→ Mean: $E[W] = 0$

$$W = \frac{Z}{\sqrt{Y/n}}$$

→ Var: $\text{Var}[W] = n/(n-2)$ for $n > 2$ (∞ for $n=1, 2$)

Assumptions:

③ Gaussian w/ unknown variance

$$CI: [M_n - \beta \sqrt{V_n}/\sqrt{n}, M_n + \beta \sqrt{V_n}/\sqrt{n}] \quad \left. \right\} \text{Student's t-distribution}$$

$$CL: 1-\alpha$$

*Find β so that $F_w(\beta) = 1 - \alpha/2 = 0.95$

④ Find Gaussian w/ unknown mean and variance

$$CI: \left[\frac{(n-1)V_n}{t_2}, \frac{(n-1)V_n}{t_1} \right] \quad \left. \right\} \text{chi-square distribution}$$

$$CL: 1-\alpha$$

*set t_1 so $F_Y(t_1) = \alpha/2$ and t_2 so $F_Y(t_2) = 1 - \alpha/2$

9.2 Significance Testing

Rejection region: $R_0: P[Y \in R_0 | H_0] = \int_{R_0} f_{Y|H_0}(y) dy = \alpha$

· If $y \in R_0$, reject the null hypothesis

· If $y \notin R_0$, fail to reject

Threshold Tests:

① 1-sided Test, Right Tail: $R_0 = \{y > t_n\}$ $P[Y > t_n | H_0] = \alpha$

② 1-sided Test, Left Tail: $R_0 = \{y < t_e\}$ $P[Y < t_e | H_0] = \alpha$

③ 2-sided Test: $R_0 = \{y > t_u\} \cup \{y < t_l\}$

$$P[Y < t_l | H_0] = P[Y > t_u | H_0] = \alpha/2$$

*Reject null hypothesis if p-value < α

① 1-sided test, Right Tail: p-value = $1 - F_{Y|H_0}(y)$

② 1-sided test, Left Tail: p-value = $F_{Y|H_0}(y)$

③ Two-sided Test: p-value = $2 \min(F_{Y|H_0}(y), 1 - F_{Y|H_0}(y))$

Four different Tests:

① One-sample Z test: observations are i.i.d. Gaussian (μ, σ^2) under H_0 .
 μ & σ^2 are known

Test: if mean has changed.

Two-sided Test: ① Compute sample mean $M_n = \frac{1}{n} \sum_{i=1}^n x_i$
② Compute Z-stat: $Z = \frac{\sqrt{n}(M_n - \mu)}{\sigma}$

③ Compute p-value: $2F_Z(-|Z|)$

④ Pval < d: Reject

Pval $\geq d$: Fail to reject

② One-sample T-test: observation i.i.d. Gaussian w/ known μ and unknown σ^2

Test: if mean has changed

Two-sided test: ① Compute sample mean

② Compute sample variance $V_n = \frac{1}{n-1} \sum_{i=1}^n (x_i - M_n)^2$

③ Compute the T-statistic $T = \frac{\sqrt{n}(M_n - \mu)}{\sqrt{V_n}} \sim T(n-1)$

④ Compute p-val: $2F_T(-|T|)$

⑤ Reject: pval < d

Fail to reject: pval $\geq d$.

③ Two-sample Z test: Collect two datasets: x_1, \dots, x_n is i.i.d. Gaussian(μ_1, σ^2) and y_1, \dots, y_m is i.i.d. Gaussian w/ the same mean $\mu_1 = \mu_2$

Means unknown variances known

Test: if means are different $\mu_1 \neq \mu_2$

Two-sided Test: ① Compute sample means

② Compute Z-stat: $\frac{M_n^{(1)} - M_n^{(2)}}{\sqrt{\sigma_1^2/n_1 + \sigma_2^2/n_2}} \sim N(0, 1)$ under H_0

③ Compute pval: $2F_Z(-|Z|)$

④ Reject: pval < d

Fail to Reject: pval $\geq d$

④ Two-sample T-test: Same as sample Z, except $\mu_1 = \mu_2$ and vars $\sigma_1^2 = \sigma_2^2$
(known μ and σ^2)

Test: if means are different $\mu_1 \neq \mu_2$

Two-sided Test: ① Compute sample means

② Compute sample vars

③ Compute: $\hat{\sigma}^2 = \frac{(n_1-1)V_n^{(1)} + (n_2-1)V_n^{(2)}}{n_1 + n_2 - 2}$

④ Compute T stat: $T = \frac{M_{n_1} - M_{n_2}}{\hat{\sigma} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$

⑤ Compute pval: $2F_T(-|T|)$

⑥ Reject pval < d, Fail pval $\geq d$