

Probability Space
 (Ω, \mathcal{F}, P)
 Sample space
 collection of events I want to assign probabilities to
 level of resolution we observe our random experiment at Axioms

F

- (i) $\Omega \in \mathcal{F}$
- (ii) If $A \in \mathcal{F}$, $A^c \in \mathcal{F}$
- (iii) If $A_i \in \mathcal{F}$, $\bigcup A_i \in \mathcal{F}$
 Hint: use of countable union to help to proof

□ $\mathcal{F} \rightarrow [0, 1]$

- (i) $0 \leq P(A) \leq 1$
 if one event occurs then other events do not
- (ii) $P(\emptyset) = 0$
- (iii) If A_1, A_2, \dots is a sequence of mutually exclusive events then
 $P(\bigcup A_i) = \sum_i P(A_i)$

Properties that follow

$$(i) P(A^c) = 1 - P(A)$$

$$\text{Proof: } P(A \cup A^c) = P(A) + P(A^c)$$

$$P(A \cup A^c) = 1 = P(A) + P(A^c)$$

$$P(A^c) = 1 - P(A)$$

$$(ii) P(\emptyset) = 0, \emptyset = \Omega^c$$

$$(iii) \text{If } A \subseteq B \text{ then } P(A) \leq P(B)$$

$$\text{Proof: } B = A \cup (A \cap B)$$

$$P(B) = P(A) + P(A \cap B) \geq P(A)$$

$$(iv) P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

$$\text{Union Bound}$$

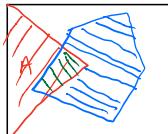
$$P(\bigcup A_i) \leq \sum_{i=1}^n P(A_i)$$

$$\text{Conditional Probability}$$

If A and B are two events, $P(B) \neq 0$, then

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

"B becomes the new universe, $B = \Omega'$ "

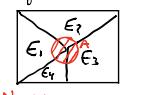


Total Probability Theorem

If $\{E_1, E_2, \dots, E_k\}$ partition Ω , then

$$P(A) = \sum_{i=1}^k P(A \cap E_i)$$

This makes more sense since you know which event which occurred



No intersections
 $\sum_{i=1}^k P(E_i) = 1$

Bayes Rule
 Want to find "ground truth" that gave this observation.
 $P(E_j|A) = \frac{P(A \cap E_j)}{P(A)} = \frac{P(E_j)P(A|E_j)}{\sum_i P(E_i)P(A|E_i)}$ prior knowledge generative model
 Independence

Two events, A_1, A_2 , are independent if $P(A_1 \cap A_2) = P(A_1)P(A_2)$

$$\Rightarrow P(A_1 | A_2) = \frac{P(A_1 \cap A_2)}{P(A_2)} = \frac{P(A_1)P(A_2)}{P(A_2)} = P(A_1)$$

Can be extended:

Events $\{A_1, A_2, \dots, A_n\}$ are independent if
 $P(A_1 \cap A_2 \cap \dots \cap A_n) = P(A_1)P(A_2) \dots P(A_n) \neq \prod_{i=1}^n P(A_i)$

$$F_{X_1, \dots, X_n}(x_1, \dots, x_n) = F_{X_1}(x_1) \dots F_{X_n}(x_n)$$

$$P[x_1 \leq x_1, \dots, x_n \leq x_n] = P[x_1 \leq x_1] \dots P[x_n \leq x_n]$$

Random Variables

Given (Ω, \mathcal{F}, P) a random variable is a function $X: \Omega \rightarrow \mathbb{R}$ such that $\forall z \in \mathbb{R}$, $\{\omega | X(\omega) \leq z\} \in \mathcal{F} \subseteq \Omega$

N^{th} Moment

N^{th} Central Moment

$$E[X^n] = \int_{-\infty}^{+\infty} x^n f_X(x) dx$$

$$E[(X - E[X])^n]$$

Cumulative Distribution Function

The CDF is defined as

$$F_X(z) = P(X \leq z) \quad \forall z \in \mathbb{R}$$

$$= P(\{\omega | X(\omega) \leq z\})$$

Properties of CDF

$$\lim_{z \rightarrow -\infty} F_X(z) = 0$$

$$\lim_{z \rightarrow +\infty} F_X(z) = 1$$

$$(2) \forall z < y \quad F_X(z) \leq F_X(y)$$

(3) F is right continuous

$$\lim_{x \rightarrow x_0^+} F_X(x) = F_X(x_0)$$

$$(4) P[x \leq X \leq y] = F_X(y) - F_X(x)$$

PMF for Discrete Random Variables

$$p_X(x) = P[X=x]$$

$$f_X(x) = \sum_{u=x} p_X(u)$$

PDF for Continuous Random Variables

$$f_X(x) = \frac{dF_X(x)}{dx} \quad \begin{array}{l} \text{measures how fast} \\ \text{we accumulate probability} \end{array}$$

$$\text{thus } F_X(x) = \int_{-\infty}^x f_X(u) du$$

Properties of PDF

$$(1) f_X(x) \geq 0 \quad \forall x$$

$$(2) \int_{-\infty}^{+\infty} f_X(x) dx = 1 = F_X(\infty)$$

$$(3) P[x \leq X \leq y] = \int_x^y f_X(u) du$$

Variance

How much a r.v. varies from its expectation

$$\text{Var}(X) = E[(X - E[X])^2]$$

$$= E[X^2 - 2E[X] + (E[X])^2]$$

$$= E[X^2] + (E[X])^2$$

$$\text{Var}(X) = \int_{-\infty}^{+\infty} g(x) f_X(x) dx = \int_{-\infty}^{+\infty} (x - E[X])^2 f_X(x) dx$$

Correlation

Correlation between X and Y

$$E[XY] = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} xy f_{X,Y}(x,y) dx dy$$

$$E[XY] = \sum_{x_i} \sum_{y_j} x_i y_j P_{X,Y}(x_i, y_j)$$

Covariance

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

$$= E[XY] - E[X]E[Y]$$

If $E[XY] = 0$, we say X is orthogonal to Y.

If $\text{Cov}(X, Y) = 0$, we say X is uncorrelated to Y.

Properties

Independence of r.v.s X, Y

$$f_{X,Y}(x, y) = f_X(x)f_Y(y) \Rightarrow \text{uncorrelatedness}$$

X, Y, Z r.v.'s

$$\text{Cov}(X, aY + bZ) = a\text{Cov}(X, Y) + b\text{Cov}(X, Z)$$

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

Expectation

$$E[X] = \begin{cases} \sum_k k P[X=k] & , X \text{ discrete} \\ \int_{-\infty}^{+\infty} x f_X(x) dx & , X \text{ continuous} \end{cases}$$

Properties

LOTUS Rule

If

$$Y = g(x)$$

then

$$E[Y] = \int_{-\infty}^{+\infty} g(x) f_X(x) dx$$

$$E[g(x)] = \int_{-\infty}^{+\infty} g(x) f_X(x) dx$$

Linearity of Expectation

$$E[\alpha X + \beta Y] = \alpha E[X] + \beta E[Y], \alpha, \beta \in \mathbb{R}$$

Preservation of Order

If

$$X \geq Y \Rightarrow E[X] \geq E[Y]$$

then

$$E[X] \geq E[Y]$$

Integration by Parts

A discrete, non-negative r.v. $X = 0, 1, 2, 3, \dots$

$$E[X] = \sum_{i=0}^{\infty} i P(X=i)$$

$$= \sum_{i=0}^{\infty} P(X \geq i) \quad \begin{array}{l} \text{Tail Probability.} \\ 1 - F_X(i) \end{array}$$

To see this, observe

$$\sum_{i=0}^{\infty} i P(X=i) = 0 \cdot P(X=0) + 1 \cdot P(X=1) + 2 \cdot P(X=2) + \dots$$

$$\sum_{i=0}^{\infty} P(X \geq i) = \begin{array}{l} \stackrel{i=0}{=} 1 \\ \stackrel{i=1}{=} \vdots \\ \vdots \\ \text{etc} \end{array} \quad \begin{array}{l} \checkmark \\ \checkmark \\ \checkmark \end{array}$$

Thus

$$E[X] = \int_0^{\infty} (1 - F_X(u)) du = \int_{-\infty}^0 F_X(u) du$$

Conditioning on Random Variables

Suppose X and Y have joint pmf $p_{X,Y}(x, y)$

The conditional pmf of X given $\{Y=y\}$ is

$$P_{X|Y}(x|y) \triangleq \Pr[X=x | Y=y]$$

$$= \frac{\Pr\{X=x \cap Y=y\}}{\Pr\{Y=y\}} = \frac{p_{X,Y}(x,y)}{p_Y(y)}$$

Similarly for continuous r.v.'s

$$f_{X|Y}(x|y) = \frac{f_{X,Y}(x,y)}{f_Y(y)} \quad \text{if } f_Y(y) \neq 0 \quad \forall y$$

Conditional Expectation

$$E[X|Y=y] = \int_{-\infty}^{+\infty} x f_{X|Y}(x|y) dx$$

$E[X|Y]$ is a random variable which takes on value $E[X|Y=y]$ w/ density $f_Y(y)$.

Since $E[X|Y]$ is r.v. can take its expectation.

$$E_Y[E_{X|Y}[X|Y]] = E[X]$$

Cauchy-Schwarz Inequality

$$|\mathbb{E}[XY]| \leq \sqrt{\mathbb{E}[X^2]}\sqrt{\mathbb{E}[Y^2]}$$

Correlation Coefficient
normalization from Cauchy-Schwarz inequality

$$\rho_{XY} = \frac{\text{Cov}(X,Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}} \quad (|\rho_{XY}| \leq 1)$$

If X_1, \dots, X_m are pairwise uncorrelated, then

$$\text{Var}(X_1 + X_2 + \dots + X_m) = \text{Var}(X_1) + \text{Var}(X_2) + \dots + \text{Var}(X_m)$$

$$\text{Var}\left(\sum_i X_i\right) = \sum_i \text{Var}(X_i)$$

Linear Minimum Mean Squared Error
Now our estimator $g(Y)$ takes the form

$$g(Y) = aY + b$$

The MSE becomes

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

$$+ \mathbb{E}[(X - aY - b)^2]$$

Want

$$\{a^*, b^*\} = \underset{a,b}{\text{argmin}} \mathbb{E}[(X - aY - b)^2]$$

Use orthogonality principle

$$W \perp g(Y) \Rightarrow X - a^*Y - b^* \perp a^*Y + b$$

Consider the following cases

$$\textcircled{1} a=0, b=1 \quad X - a^*Y - b^* + 1 \Rightarrow \mathbb{E}[X - a^*Y - b^*] = 0$$

$$\mathbb{E}[X] - \mathbb{E}[a^*Y] - \mathbb{E}[b^*] = 0 \quad \text{W is UNBIASED}$$

$$\Rightarrow b^* = \mathbb{E}[X] - a^*\mathbb{E}[Y]$$

$$\textcircled{2} a=1, b=0$$

$$X - a^*Y - b^* \perp Y \Rightarrow \mathbb{E}[(X - a^*Y - b^*)Y] = 0$$

$$\mathbb{E}[(X - a^*Y - b^*)Y]$$

$$= \mathbb{E}[(X - a^*Y - \mathbb{E}[Y]) - (\mathbb{E}[X] - a^*\mathbb{E}[Y])]Y$$

$$\stackrel{\text{cancel } a^*Y}{=} \mathbb{E}[(X - \mathbb{E}[X]) - (\mathbb{E}[Y] - \mathbb{E}[Y])]Y$$

$$\Rightarrow \text{Cov}(X - \mathbb{E}[X], Y - \mathbb{E}[Y]) = 0$$

$$\text{Cov}(X, Y) = a^* \text{Cov}(X, Y) \quad \text{all } X_i \text{ pairwise orthogonal}$$

$$a^* = \text{Cov}(X, Y) / \text{Var}(Y) \quad \text{Cov}(X, Y)$$

Aside

$$\text{Cov}(X, aY + bZ) = a\text{Cov}(X, Y) + b\text{Cov}(X, Z)$$

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

$$\textcircled{3} \quad \hat{X}_{\text{MSE}} = a^*Y + b^*$$

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

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

and the MSE is

$$\text{MSE} = \mathbb{E}[W^2] = \text{Cov}(W, W)$$

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

$$= \text{Cov}(X, X) + 2a^* \text{Cov}(X, Y) - 2a^* \text{Cov}(X, Y)$$

$$= \text{Var}(X) + \frac{(\text{Cov}(X, Y))^2}{\text{Var}(Y)} \text{Var}(Y) - 2 \frac{(\text{Cov}(X, Y)) \text{Cov}(X, Y)}{\text{Var}(Y)}$$

$$= \text{Var}(X) - \frac{(\text{Cov}(X, Y))^2}{\text{Var}(Y)}$$

$$W = X - \hat{X}_{\text{MSE}}$$

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

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

Know

$W \perp Y$ are Jointly Gaussian

$$\mathbb{E}[WY] = 0 \quad (\text{orthogonality})$$

so $\text{Cov}(W, Y) = 0 \Rightarrow$ independence by property 2

Need W is error given by best MSE

$$\mathbb{E}[Wg(Y)] = 0 \quad \text{if } g(Y) \Rightarrow$$

Trivially true since W independent of Y .

$$\mathbb{E}[Wg(Y)] = \mathbb{E}[W] \mathbb{E}[g(Y)] = 0$$

Proof of ④

Hilbert Space

A Hilbert space is a vector space that

- ① Has an inner product $\langle \cdot, \cdot \rangle$ defined
- ② Is complete w.r.t respect to the norm induced by $\langle \cdot, \cdot \rangle$

$\|x\| = \sqrt{\langle x, x \rangle}$ complete in the sense every Cauchy sequence converges

Hilbert Space of Random Variables

All random variables w.r.t finite 2nd moments form a Hilbert space w.r.t inner product $\langle x, y \rangle \triangleq \mathbb{E}[xy]$

$$\langle x, y \rangle \triangleq \mathbb{E}[xy]$$

Inner Product \Rightarrow norm/length \Rightarrow distance metric

$$\text{length/norm of } X = \sqrt{\langle X, X \rangle} = \sqrt{\mathbb{E}[X^2]} \quad \text{Cauchy-Schwarz} \quad \text{this} \leq 1$$

$$\text{Angle b/w } X, Y; \theta_{XY}: \cos(\theta_{XY}) = \frac{\langle X, Y \rangle}{\|X\| \|Y\|} = \frac{\mathbb{E}[XY]}{\sqrt{\mathbb{E}[X^2] \mathbb{E}[Y^2]}}$$

Projection of X onto Y

$$\hat{X}_Y(X) = \frac{\langle X, Y \rangle}{\|Y\|^2} Y = \frac{\mathbb{E}[XY]}{\mathbb{E}[Y^2]} Y \quad \text{a random variable}$$

Orthogonality

$$\langle X, X \rangle \triangleq \mathbb{E}[XX] = 0$$

$$\hat{X}_Y(X) = 0$$

Pythagorean

$$\mathbb{E}\left[\left(\sum_i X_i\right)^2\right] = \sum_i \mathbb{E}[X_i^2] \quad \text{pairwise orthogonal}$$

Jointly Gaussian Random Variables

If

$$X \sim N(\mu, \sigma^2)$$

then

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

Note: K is our covariance matrix.

Here, $K = \begin{pmatrix} \text{Cov}(X, X) & \text{Cov}(X, Y) \\ \text{Cov}(Y, X) & \text{Cov}(Y, Y) \end{pmatrix}$

$$= \mathbb{E}\left[\begin{pmatrix} X & Y \\ Y & X \end{pmatrix} \begin{pmatrix} X & Y \\ Y & X \end{pmatrix}^\top\right]$$

$$= \frac{\sigma^2}{\sqrt{2\pi\sigma^2}} \text{diag}(K)$$

Definition: X, Y are jointly Gaussian if $aX + bY$ is Gaussian $\forall a, b \in \mathbb{R}$.

Definition: Joint pdf of such a distribution is

$$f_{X,Y}(x,y) = \frac{1}{2\pi\sqrt{|K|}} \exp\left(-\frac{1}{2} \begin{pmatrix} x - \mu_x \\ y - \mu_y \end{pmatrix}^\top K^{-1} \begin{pmatrix} x - \mu_x \\ y - \mu_y \end{pmatrix}\right)$$

Properties of Jointly Gaussian Random Variables

$$\textcircled{1} \quad \text{J.G.} \Rightarrow \text{M.G.}$$

$\textcircled{2}$ Uncorrelated JG r.v.s are independent **special!**
independence \Rightarrow uncorrelated
not in general

$\textcircled{3}$ If X, Y jointly Gaussian, then

$$a_1X + b_1Y + c_1, \quad a_2X + b_2Y + c_2 \quad \text{are jointly Gaussian}$$

$\textcircled{4}$

For J.G. X, Y

$$\hat{X}_{\text{MSE}} = \mathbb{E}[X|Y] = \hat{X}_{\text{MSE}} = \mathbb{E}[X] + \frac{\text{Cov}(X, Y)}{\text{Var}(Y)} (Y - \mathbb{E}[Y])$$

$$\text{MSE} = \text{Var}(X) (1 - \frac{\text{Cov}(X, Y)}{\text{Var}(Y)}) = \text{Var}(X) - \frac{\text{Cov}^2(X, Y)}{\text{Var}(Y)}$$

$$\textcircled{5} \quad \text{For J.G. } X, Y \quad \mathbb{E}[X|Y=y] = \mathbb{E}[X] + \frac{\text{Cov}(X, Y)}{\text{Var}(Y)} (y - \mathbb{E}[Y])$$

$$\mathbb{E}[X|Y] = \mathbb{E}[X] + \frac{\text{Cov}(X, Y)}{\text{Var}(Y)} (Y - \mathbb{E}[Y])$$

$$\text{Var}(X|Y=y) = \text{Var}(X) - \frac{\text{Cov}^2(X, Y)}{\text{Var}(Y)}$$

$$\text{Var}(X|Y) = \text{Var}(X) - \frac{\text{Cov}^2(X, Y)}{\text{Var}(Y)}$$

$\textcircled{6}$ The conditional distribution of X given $Y=y$ is Gaussian

$$f_{X|Y}(x|y) = \frac{1}{\sqrt{2\pi\text{Var}(X|Y)}} e^{-\frac{(x-\mu_{X|Y})^2}{2\text{Var}(X|Y)}}$$

Mean Squared Error (MSE)

Have r.v. X which can't be directly observed.

Want to estimate \hat{X} .

Note: we are in a Hilbert space of random variables

The estimation error

$$W = X - \hat{X}$$

The MSE of \hat{X} is

$$\mathbb{E}[(X - \hat{X})^2] = \mathbb{E}[W^2] = \langle W, W \rangle = \|W\|^2 = (\text{d}(X, \hat{X}))^2$$

Minimum Mean Squared Error

$$\hat{X}_{\text{MSE}} = \underset{X}{\operatorname{argmin}} \mathbb{E}[(X - \hat{X})^2]$$

MSE of X using a constant

$$\hat{X}_{\text{MSE}} = a^*$$

where

$$a^* = \underset{a \in \mathbb{R}}{\operatorname{argmin}} \mathbb{E}[X - a]^2$$

it's $\mathbb{E}[X]$!

$$\text{MSE} = \mathbb{E}[(X - a)^2]$$

$$= \mathbb{E}[(X - \mathbb{E}[X])^2] + \mathbb{E}[\mathbb{E}[X] - a]^2$$

$$= \mathbb{E}[(X - \mathbb{E}[X])^2] + 2\mathbb{E}[(X - \mathbb{E}[X])(\mathbb{E}[X] - a)] + \mathbb{E}[(\mathbb{E}[X] - a)^2]$$

Want to minimize

$$\mathbb{E}[(X - \mathbb{E}[X])^2] + \mathbb{E}[\mathbb{E}[X] - a]^2$$

$$\rightarrow \text{Var}(X) + (\mathbb{E}[X] - a)^2$$

$$\text{Take } a = \mathbb{E}[X]$$

MSE of X given $Y=y$

Similar proof

$$\hat{X}_{\text{MSE}} = \mathbb{E}[X|Y]$$

$$\text{NSE}(\hat{X}_{\text{MSE}}) = \mathbb{E}[\text{Var}(X|Y)]$$

The error

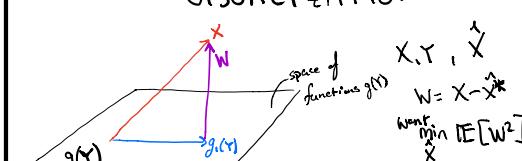
$$W = X - \mathbb{E}[X|Y]$$

has

$$\mathbb{E}[W] = \mathbb{E}[X - \mathbb{E}[X|Y]] = \mathbb{E}[X] - \mathbb{E}[X] = 0$$

we call this type of estimator **UNBIASED**

VISUALIZATION



- Orthogonality principle:

$$W + g(Y) \perp g(Y) \quad \text{necessary AND sufficient condition}$$

W induced by X^* (best estimator)

$$\text{claim } X - \mathbb{E}[X|Y] \perp g(Y) \perp g \quad \text{inner product } \langle X - \mathbb{E}[X|Y], g(Y) \rangle$$

$$\text{LHS} \quad \mathbb{E}_Y \left[X g(Y) \right] - \mathbb{E}_Y \left[(\mathbb{E}_Y [X|Y]) g(Y) \right]$$

$$= \mathbb{E}_Y \left[\mathbb{E}_{X|Y} [X g(Y)|Y] \right] - \mathbb{E}_Y \left[\mathbb{E}_{X|Y} [\mathbb{E}_Y [X|Y]|Y] g(Y) \right]$$

$$= \mathbb{E}_Y [g(Y) \mathbb{E}_{X|Y} [X|Y]] - \mathbb{E}_Y [g(Y) \mathbb{E}_{X|Y} [X|Y]] = 0 \quad \checkmark$$

MSE Using Pythagorean Theory

$$\|W\|^2 = \langle W, W \rangle \text{ in Hilbert Space}$$

$$= \mathbb{E}[W^2] = \mathbb{E}[X^2] - \mathbb{E}[X^2]$$

$$= \text{Var}(X) + (\mathbb{E}[X])^2 - (\text{Var}(X) + (\mathbb{E}[X])^2)$$

$$= \text{Var}(X) - \text{Var}(\hat{X}_{\text{MSE}})$$

$$\mathbb{E}[X] = \mathbb{E}[\hat{X}_{\text{MSE}}] = \text{Var}(X) - \text{Var}(\hat{X}_{\text{MSE}})$$

- \hat{X}_{MSE} is a constant Thus $X \sim N(\mu_W + \hat{X}_{\text{MSE}}, \sigma^2)$

- W maintains its distribution

