

Decision Theory

States

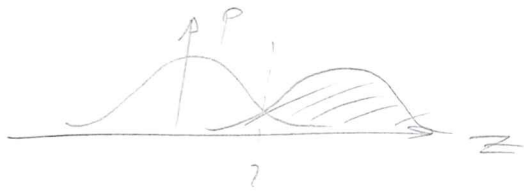
measurements

s_1

$P(z|s_1)$ - "healthy"

s_2

$P(z|s_2)$ - "ill, sick"



Task

decide d_1 (guessed s_1) or d_2 (s_2) based on z .



$$\Lambda(z) = \frac{P(z|s_2)}{P(z|s_1)} \equiv \text{likelihood ratio}$$

Maximum likelihood test

d_1 if $P(z|s_1) > P(z|s_2)$
 d_2 else

$$\Lambda(z) \begin{matrix} \geq 1 \\ \rightarrow d_1 \\ < 1 \\ \rightarrow d_2 \end{matrix}$$

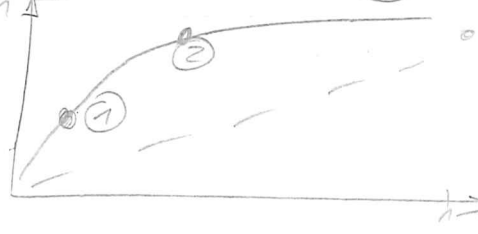
Types of error

$P(d_2|s_1) \equiv$ Type-I-error ("false positive", "false alarm")

$P(d_1|s_2) \equiv$ II ("false negative")

Receiver-operating-characteristic (ROC)

$P(d_2|s_2)$
detection level



① type-I-error expensive (eg. bus enforcement)

② type-II error expensive (medical diagnosis)

Bayes risk criterion

- Cost matrix

C_{11} = cost of d_1 if S_1 true

C_{12} = " " " " d_1 if S_2 true

- base rates $P(S_1), P(S_2)$

- expected cost

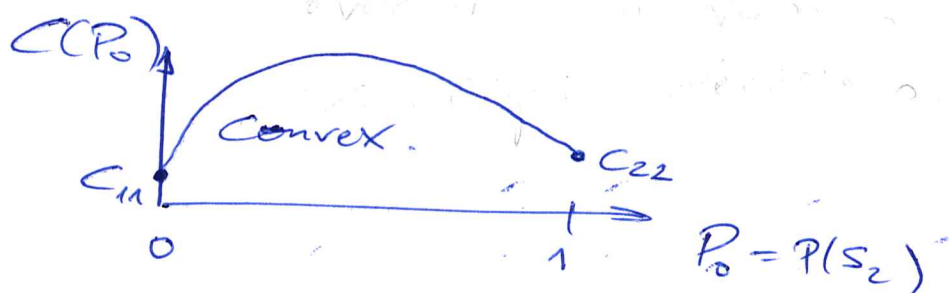
$$\langle C \rangle = C_{11} P(d_1 | S_1) \cdot P(S_1) +$$

$$C_{21} P(d_2 | S_1) P(S_1) + \dots$$

- optimal test for $\overrightarrow{\text{likelihood}}^{\text{min}}$ ratio

$$\frac{P(z|S_2)}{P(z|S_1)} = \Delta(z) \geq \frac{d_2 (C_{21} - C_{11}) P(S_1)}{d_1 (C_{12} - C_{22}) P(S_2)}$$

- expected cost for optimal test



Generalizations

- vector observations

- multiple decisions d_1, d_2, d_3, \dots
for S_1, S_2, S_3, \dots

• Q_1, Q_2 : (\equiv critical base rates)

$$C_{12} Q_1 = C(Q_1) + \beta$$

$$C_{21}(1-Q_2) = C(Q_2) + \beta$$

expected cost for...

• no observation

• one observation

$$T_1 = \frac{Q_1}{1-Q_1} \quad T_2 = \frac{Q_2}{1-Q_2}$$

(\equiv corresponding likelihood ratios)

$$Z_1 = \{1_0 < T_1\}$$

$$Z_2 = \{1_0 > T_2\}$$

• Next step.

$$\begin{aligned} \Lambda_1(z) &= \frac{P_1}{1-P_1}, \quad P_1 = P(S_2 | Z_1) \\ &= \frac{P(S_2 | Z_1)}{P(S_1 | Z_1)} \end{aligned}$$

Assumptions: observations indep.
and identically distributed

\Rightarrow same statistics at each step, but prior updated.

$$Z_1 = \{z_k: \Lambda_k(z_k) < T_1\}$$

$$Z_2 = \{z_k: \Lambda_k(z_k) > T_2\}$$

Decision theory

$n \geq 2$ ~~two~~ alternative measurement models

↓
decide which one more likely

Statistical test

- Working hypothesis H_1
(e.g. smoking promotes cancer)
→ no measurement model

- null hypothesis $H_0 = \neg H_1$
measurement model

$p(z)$, $z \equiv$ measurement variable
if H_0 true.

Statistical tests

H_0 : null hypothesis

Statistical model of reality

$P(z | H_0)$, $P(z | H_1)$

Statistical test = decision problem

d_1 : accept H_0 : $P(z | H_0)$

d_2 : reject H_0 : no model

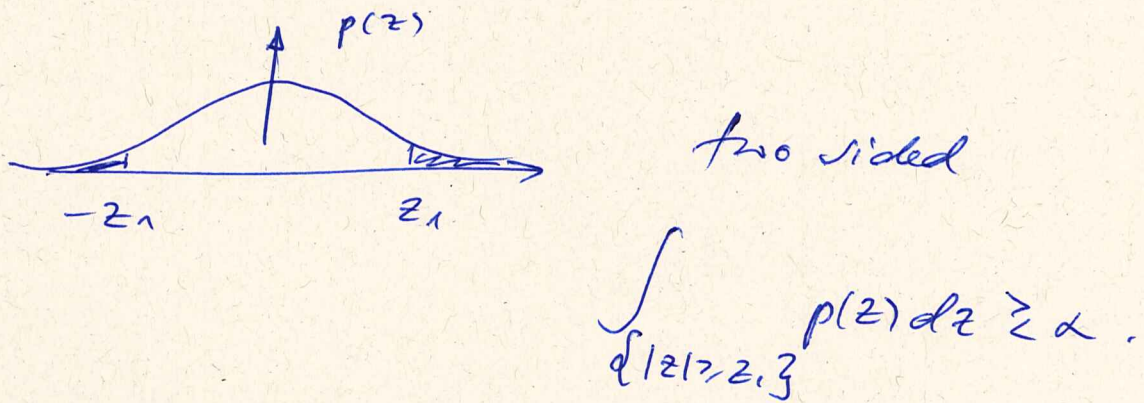
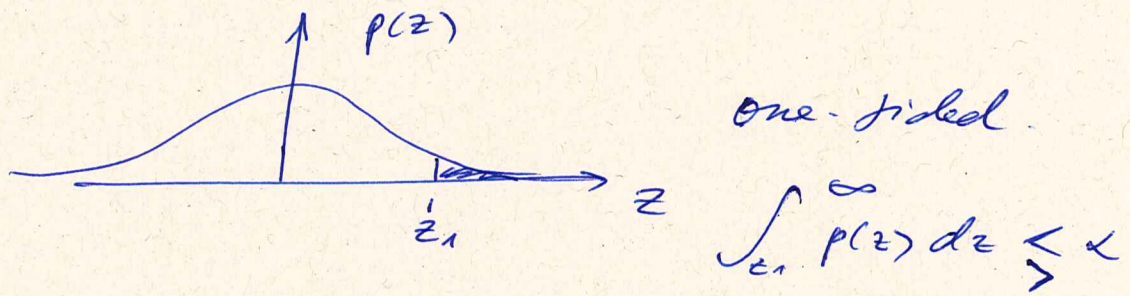
decision rule

observation z_0
 $P(z \geq z_0 | H_0) \stackrel{d_1}{>} \alpha \stackrel{d_2}{<}$, $\alpha =$ significance level

- "reject H_0 if z_0 too unlikely"
- $\alpha =$ probability of false positives.
- "look-elsewhere fallacy"

Simultaneous testing of n null hypotheses:
 H_1, \dots, H_n with fixed α
even if $H_j, j=1, \dots, n$ true,
on average $\alpha \cdot n$ hypotheses will be rejected.

• One-sided / two-sided tests.



Which depends on H_1 !

Example (fair coin)

$$H_0: p_{\text{tail}} = p_{\text{head}} = \frac{1}{2}$$

Measurement

k -times 'tail' after n -times measurement

measurement model

$$p(k) = \frac{\binom{n}{k}}{2^n} \equiv \text{Binomial distribution}$$

one-sided test

$$\sum_{k=0}^{k_1} p(k) < \alpha$$

H_1 : bias towards tail

two-sided test

$$\sum_{k=k_1}^n + \sum_{k=1}^{k_1} p(k) < \alpha$$

H_1 : unknown bias

Student's t-test

1908: William Gosset at "Guinness"
(pen name "a student")

$$H_0: \langle X \rangle = \mu$$

measurement model:

$$z = X \sim \mathcal{N}(\mu, \sigma^2).$$

with σ^2 unknown.

measurement

$$z_1, \dots, z_n.$$

test

$$t = \frac{\bar{z} - \mu}{\hat{\sigma} / \sqrt{n}}$$

μ assumed to be known (by H_0)

$\hat{\sigma}^2$ = estimated variance

- ~~large~~ independent of μ, σ^2 .
- strategy one: Stochastic simulations.
- — — — two: analytic formula
(with simplifying assumption)
 - numerator & denominator indep.
 - denominator $\sim \chi^2(u-1)$

Generalization

• two populations: $H_0: \mu_1 = \mu_2$

unpaired: e.g. group of males, group of females

paired: e.g. group of patients before/after treatment

Def: χ^2 -distribution

$$y_1, \dots, y_n \sim W(0, 1).$$

\Rightarrow

$$S = y_1^2 + \dots + y_n^2 \sim \chi^2(n)$$

mean: n

variance: $2n$.

Application to $\hat{\sigma}^2$:

$$z_1, \dots, z_n \sim W(\mu, \sigma^2)$$

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_{j=1}^n (z_j - \bar{z})^2 \sim \sigma^2 \cdot \underbrace{\chi^2(n-1)}_{n-1}$$

mean: 1 .

Proof: $(z_1, \dots, z_n) \in \mathbb{R}^n \equiv$ spherically distributed
multivariate normal distribution.

↓ linear
proj.
operator.

$$(z_1 - \bar{z}, \dots, z_n - \bar{z}) \equiv \text{---} \text{---}$$

$$\uparrow W(0, \frac{n-1}{n} \sigma^2)$$

mean: 0

on $(n-1)$ -dim'l
subspace.

$$\text{Variance: } \left(\frac{n-1}{n}\right)^2 \sigma^2 + (n-1) \left(\frac{1}{n}\right)^2 \sigma^2 = \frac{n-1}{n} \sigma^2.$$

Note: $\sum_{j=1}^n (x_j - \mu)^2 = \sum_{j=1}^n (x_j - \bar{x})^2 + n(\bar{x} - \mu)^2$ easy proof

Fisher's exact test:

testing categorical data for correlation

• Categorical properties

- A; e.g. smoking

B; e.g. cancer.

• Contingency table (absolute frequencies)

	A	\bar{A}	
B	N_{AB}	$N_{\bar{A}B}$	N_B ← margin totals.
\bar{B}	$N_{A\bar{B}}$	$N_{\bar{A}\bar{B}}$	$N_{\bar{B}}$
	N_A	$N_{\bar{A}}$	N

• $H_0: p(B|A) = p(B)$

• Simplification: Suppose N_A, N_B, N fixed.

⇒ Single degree of freedom, e.g. N_{AB} .

$$P(N_{AB}) = \frac{\binom{N_B}{N_{AB}} \binom{N_{\bar{B}}}{N_A - N_{AB}}}{\binom{N}{N_A}} \equiv \text{hyper-geometric distribution}$$

(given N elements.
in total out of which
K are positive; what is
prob. to have k positive
in selection of n)

Estimation theory

- Measurement model
 $p(z|x)$, $z \equiv$ measurement
 $x \equiv$ unknown parameter
- Max likelihood estimate
 \hat{x} that maximizes $p(z|x)$
- Later: Bayesian prior $p(x)$.

Warm-up example

How to combine using measurements?

- Methods 1 and 2 to measure x

$$P(z_1|x) = N_{z_1}(x, \sigma_1^2) \leftarrow \text{variance } \sigma_1^2$$

$$P(z_2|x) = N_{z_2}(x, \sigma_2^2) \leftarrow \text{--- } \sigma_2^2$$

- Max-likelihood estimate for x

$$P(z_1, z_2 | \hat{x}) = P(z_1 | \hat{x}) \cdot P(z_2 | \hat{x}) \rightarrow \max$$

\Rightarrow

$$\hat{x} = \frac{\sigma_2^2}{\sigma_1^2 + \sigma_2^2} z_1 + \frac{\sigma_1^2}{\sigma_1^2 + \sigma_2^2} z_2$$

\equiv weighted mean

NB: z_1, \dots, z_n with $\sigma_1^2 = \dots = \sigma_n^2$

\rightarrow arithmetic mean

Next: A-priori knowledge on $p(x)$
 \equiv Bayesian prior

Bayes formula

$$P(x|z) = \frac{P(z|x)P(x)}{P(z)}$$

$P(x)$ = Bayesian prior.

$P(z|x)$ = measurement model

$$P(z) = \int dx P(z|x)P(x)$$

$P(x|z)$ = a-posteriori prob. distrib.

Example

$$P(x) = N(0, \sigma^2)$$

$$P(z|x) = N(x, \theta^2) \Rightarrow P(z) = N(0, \sigma^2 + \theta^2)$$

\Rightarrow

$$P(x|z) = N(\hat{x}, p^2)$$

$$\hat{x} = \frac{1}{1 + \theta^2/\sigma^2} z, \quad p^2 = \frac{1}{\frac{1}{\sigma^2} + \frac{1}{\theta^2}}$$

\hat{x} = weighted mean:

$$\hat{x} = \frac{\theta^2}{\sigma^2 + \theta^2} \cdot 0 + \frac{\sigma^2}{\sigma^2 + \theta^2} z$$

Sequential measurements.

Use previous a-posteriori distribution as new prior.

Example: Infotaxis, Kalman filter.

Curve fitting = estimation problem

$$x_i = f_i(\alpha), \quad i=1, \dots, k$$

data points. physical model

$\alpha \equiv$ unknown parameters

- noisy measurements

$$z_i = x_i + \epsilon_i \sim W(x_i, \sigma_i^2) \quad \sigma_i^2 = \sigma^2$$

- likelihood for observation

$$P(\{z_i\} | \alpha) \rightarrow \text{max for } \alpha = \hat{\alpha}$$

$$P(\{z_i\} | \alpha) = \prod_{i=1}^k N_{z_i}(f_i(\alpha), \sigma^2)$$

$$\sim \exp \left[- \frac{\sum_{i=1}^k (f_i(\alpha) - z_i)^2}{2\sigma^2} \right]$$

$$\Rightarrow \alpha \text{ maximizes } \sum_{i=1}^k (f_i(\alpha) - z_i)^2$$

\Rightarrow least-square fit.

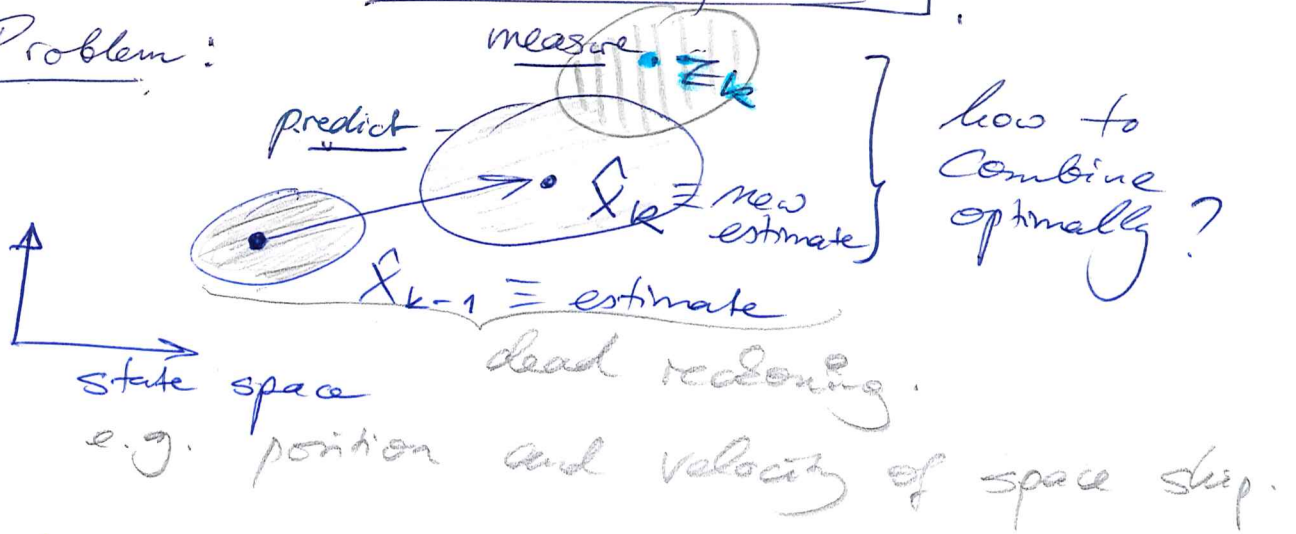
generalizations:

- different $\sigma^2 \Rightarrow$ weighted fit

- non-gaussian noise
maximize $P(\{z_i\} | \alpha)$

Kalman filter

Problem:



Application:

- Apollo.
- Missile guidance (Tomahawk)
- GPS.

Kalman filter

Notation

$\hat{x}_{k-1|k-1}$ \equiv estimate of state at time t_{k-1}

$P_{k-1|k-1}$ \equiv covariance matrix of estimate

↓ predict

$\hat{x}_{k|k-1}$ \equiv a-priori estimate for time t_k .

$P_{k|k-1}$ \equiv a-priori-covariance

↓ measure
 z_k

update (weighted average)

$\hat{x}_{k|k}$ \equiv a-posteriori estimate

$P_{k|k}$ \equiv a-posteriori covariance

Predict

$$\hat{X}_{k|k-1} = F_k X_{k|k-1} \equiv \text{a-priori}$$

$$P_{k|k-1} = F_k P_{k-1|k-1} F_k^T + Q_k$$

Measure

$$z_k$$

$$y_k = z_k - H_k \hat{X}_{k|k-1} \equiv \text{pre-filter residual (innovation)}$$

$$S_k = R + H_k P_{k|k-1} H_k^T \equiv \text{covariance of } y_k$$

Update:

$$\hat{X}_{k|k} = \hat{X}_{k|k-1} + K_k y_k \equiv \text{weighted average}$$

$$P_{k|k} = (I - K_k H_k) P_{k|k-1} (I - K_k H_k)^T + K_k R_k K_k^T$$

\equiv Covariance

\rightarrow a-posteriori.

Optimal Kalman gain.

$$K_k = \underset{n \times m}{P_{k|k-1}} \underset{n \times m}{H_k^T} \underset{m \times m}{S_k^{-1}}$$

\approx "ratio of covariances"

Theorem

The optimal Kalman gain

minimizes $E \|x_k - \hat{x}_{k|k}\|^2$.

N.B. For the optimal gain $P_{k|k} = (I - K_k H_k) P_{k|k-1}$.

N.B. Only update variables that have been measured (or are correlated to a measurement).

Process Model

$$X_k = F_k X_{k-1} + W_k$$

$F_k \equiv$ time propagator

$W_k \equiv$ process noise

$W_k \sim \mathcal{N}(0, Q_k)$, i.e. mean 0
covariance Q_k .

Measurement model

$$z_k = H_k X_k + v_k$$

$H_k \equiv$ measurement matrix

N.B. w.l.o.g. $H_k = \begin{pmatrix} 1 & & & & \\ & \dots & & & \\ & & 1 & & \\ & & & \dots & \\ & & & & 0 \end{pmatrix}$

i.e. only $x_{k,1}, \dots, x_{k,u}$, $u \leq n$
are measured

$v_k \equiv$ measurement noise

$v_k \sim \mathcal{N}(0, R_k)$.

Proof

$$E \|x_k - \hat{x}_{k|k}\|^2 = \text{tr } P_{k|k}$$

$$(*) \quad 0 \stackrel{!}{=} \frac{\partial \text{tr } P_{k|k}}{\partial K_k}$$

We need rules of matrix gradients:

$$\bullet \frac{\partial \text{tr } AB}{\partial A} = B^T$$

because

$$\left(\frac{\partial \text{tr } AB}{\partial A} \right)_{ij} = \frac{\partial \text{tr } AB}{\partial A_{ij}} = \frac{\partial \sum_k (AB)_{kk}}{\partial A_{ij}} = \frac{\partial \sum_r A_{ir} B_{rk}}{\partial A_{ij}} = B_{ji}$$

$$\bullet \frac{\partial \text{tr } B A^T}{\partial A} = B$$

Now

$$\begin{aligned} \text{tr } P_{k|k} &= P_{k|k-1} - K_k H_k^T P_{k|k-1} - P_{k|k-1} (K_k H_k^T)^T \\ &\quad + \underbrace{K_k H_k^T P_{k|k-1} (K_k H_k^T)^T + K_k R_k K_k^T}_{K_k S_k K_k^T} \end{aligned}$$

$$\begin{aligned} \frac{\partial \text{tr } P_{k|k}}{\partial K_k} &= -P_{k|k-1}^T H_k^T - P_{k|k-1} H_k^T + K_k S_k^T \\ &\quad + K_k S_k \\ &= -2 P_{k|k-1} H_k^T + 2 K_k S_k \end{aligned}$$

From (*), $K_k = -P_{k|k-1} H_k^T S_k^{-1}$

q.e.d.

Kalman filter example

Velocity $T\dot{v} = -v + \xi$ } discrete time
position $\dot{x} = v$ } $v_k = v_{k-1} - v_{k-1} \cdot \frac{\Delta t}{T}$
 $+ \underbrace{\sqrt{\frac{2D\Delta t}{T^2}}}_{W_k} N_k$
Measurement $x_k = x_{k-1} + v_{k-1} \cdot \Delta t$

$$z_k = x_k + \theta_k, \quad x_k = x(t_k)$$

→ Matrix formalism

$$X_k = \begin{pmatrix} v_k \\ x_k \end{pmatrix}$$

$$X_k = F_k X_{k-1} + W_k$$

$$F_k = \begin{pmatrix} 1 - \frac{\Delta t}{T} & 0 \\ \Delta t & 0 \end{pmatrix}, \quad W_k \sim \mathcal{N}(0, Q_k)$$
$$Q_k = \begin{pmatrix} 2D\Delta t/T^2 \\ 0 \end{pmatrix}$$

$$z_k = H_k x_k + v_k$$

$$H_k = (0 \ 1)$$

$$v_k \sim \mathcal{N}(0, R_k)$$

$$R_k = D_z$$

Exercise

- Simulate X_k, z_k , plot.
- Compute $\hat{x}_{k|k-1}, y_k, S_k, \hat{x}_{k|k}, P_{k|k}, K_k$.
- plot $\hat{x}_{k|k}$: effect of D_z ?
- measure also velocity