

Chapter 2: Fundamentals of prediction aka Decision Theory

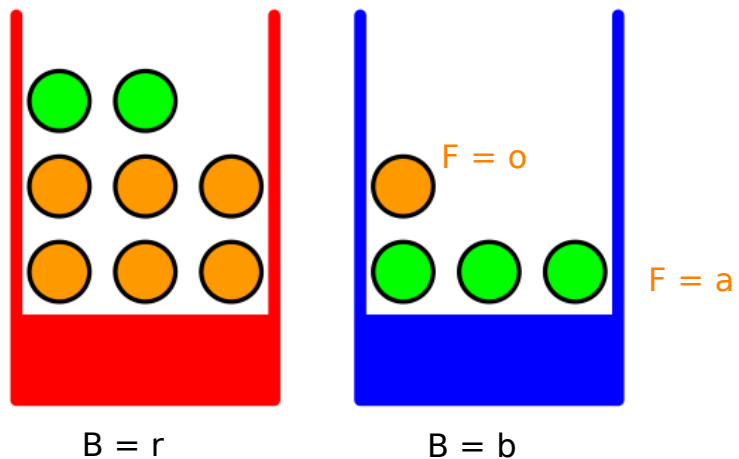
Patterns, Predictions and Actions by Hardt and Recht

Chapter 1: Pattern recognition; by Bishop

1. Probability Basics

We will introduce the basic concepts of probability theory by considering a simple example. Imagine we have two boxes, one red and one blue, and in the red box we have 2 apples and 6 oranges, and in the blue box we have 3 apples and 1 orange. This is illustrated in Figure 1.9. Now suppose we randomly pick one of the boxes and from that box we randomly select an item of fruit, and having observed which sort of fruit it is we replace it in the box from which it came. We could imagine repeating this process many times. Let us suppose that in so doing we pick the red box 40% of the time and we pick the blue box 60% of the time, and that when we remove an item of fruit from a box we are equally likely to select any of the pieces of fruit in the box.

Figure 1.9 We use a simple example of two coloured boxes each containing fruit (apples shown in green and oranges shown in orange) to introduce the basic ideas of probability.



$$p(B = r) = 4/10 \text{ and } p(B = b) = 6/10$$

$$p(F = a | B = r) = 2/8 \text{ and } p(F = o | B = r) = 6/8$$

$$p(F = a | B = b) = 3/4 \text{ and } p(F = o | B = b) = 1/4$$

Joint prob

$$p(F=a, B=r) = p(F=a | B=r) p(B=r) \leftarrow \text{product rule}$$

Marginal prob

$$p(F=a) = \sum_{B \in \{r, b\}} p(F=a, B=B) \leftarrow \text{sum rule}$$

$$= p(F=a, B=r) + p(F=a, B=b)$$

$$= p(F=a | B=r) p(B=r) + p(F=a | B=b) p(B=b)$$

$$P(F=a|B=r) \xrightarrow{\text{Bayes rule}} P(B=r|F=a)$$

$$P(B=r|F=a) = \frac{P(B=r, F=a)}{P(F=a)}$$

$$P(B=r|F=a) = \frac{P(F=a|B=r)P(B=r)}{P(F=a)}$$

The Rules of Probability

sum rule

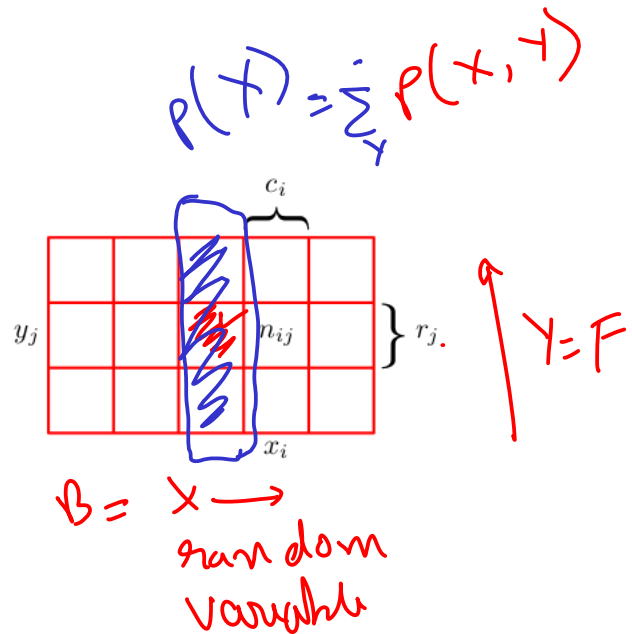
$$p(X) = \sum_Y p(X, Y) \quad (1.10)$$

Bayes rule

product rule

$$p(X, Y) = p(Y|X)p(X). \quad (1.11)$$

Figure 1.10 We can derive the sum and product rules of probability by considering two random variables, X , which takes the values $\{x_i\}$ where $i = 1, \dots, M$, and Y , which takes the values $\{y_j\}$ where $j = 1, \dots, L$. In this illustration we have $M = 5$ and $L = 3$. If we consider a total number N of instances of these variables, then we denote the number of instances where $X = x_i$ and $Y = y_j$ by n_{ij} , which is the number of points in the corresponding cell of the array. The number of points in column i , corresponding to $X = x_i$, is denoted by c_i , and the number of points in row j , corresponding to $Y = y_j$, is denoted by r_j .



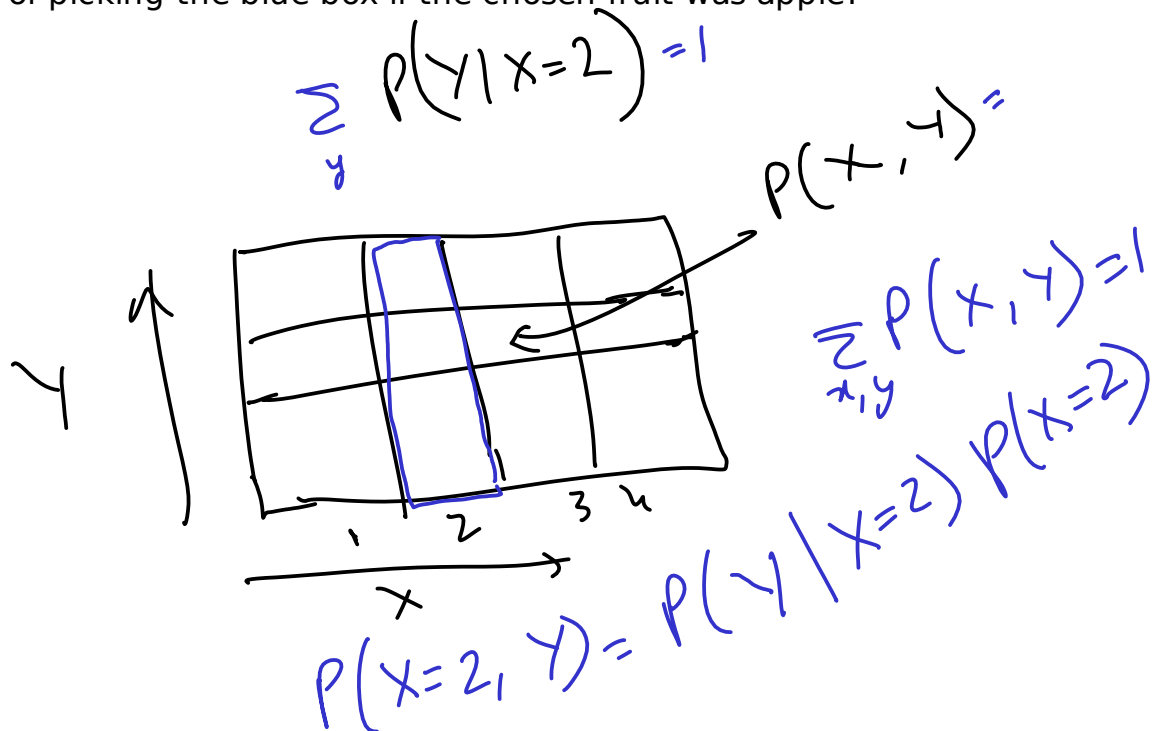
What is Bayes rule?

Homework 5

Find the probability of picking the blue box if the chosen fruit was apple?

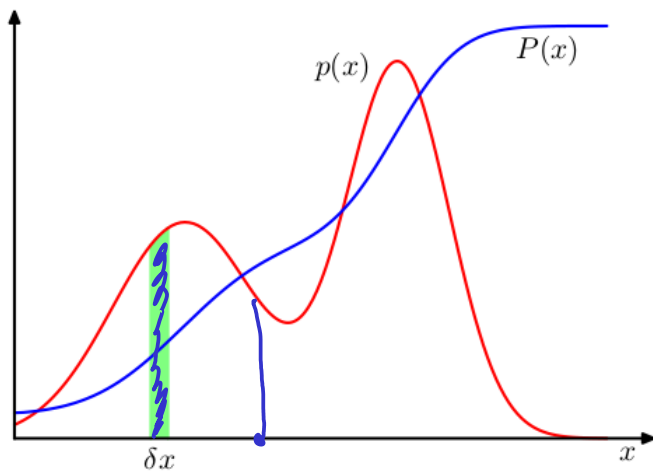
Problem 1

Product rule



Probabilities of continuous variables

Figure 1.12 The concept of probability for discrete variables can be extended to that of a probability density $p(x)$ over a continuous variable x and is such that the probability of x lying in the interval $(x, x + \delta x)$ is given by $p(x)\delta x$ for $\delta x \rightarrow 0$. The probability density can be expressed as the derivative of a cumulative distribution function $P(x)$.



$$P(x = 1.7349...) = 0 \quad x = \text{cont. RV}$$

$$P(x \in [1.73, 1.73 + 0.01]) = \text{some value}$$

$$f(x) = \text{Prob density function} = \lim_{\delta x \rightarrow 0} \frac{P(x \in [x, x + \delta x])}{\delta x}$$

$$P(x \in [x, x + \delta x]) = \int_x^{x + \delta x} f(x) dx$$

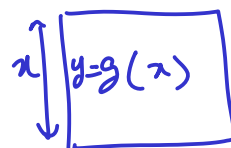
$\delta x = 0.01$

$$0 \leq f(x)$$

$$f(x) \text{ can be } > 1$$

Transformation of random variables

$$\begin{aligned} p_y(y) &= p_x(x) \left| \frac{dx}{dy} \right| \\ &= p_x(g(y)) |g'(y)|. \end{aligned}$$



X is side of square

$$P_x(x=x) = f_x(x) = \text{Prob. density over } x$$

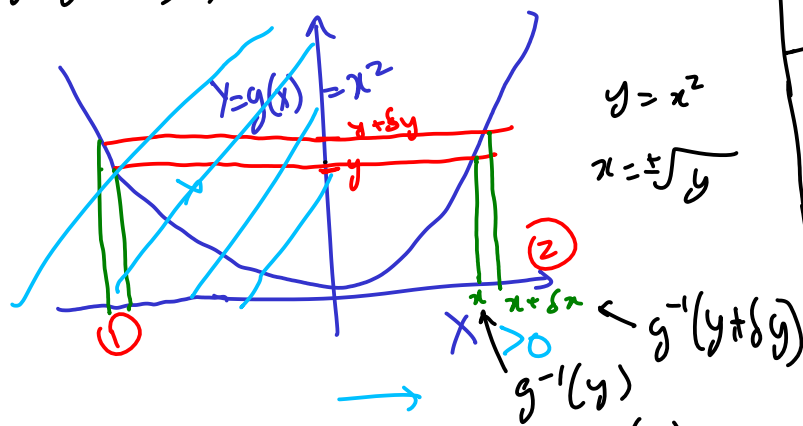
$$y = g(x) = x^2$$

↑
is another RV

Transformation

Find the PDF of Y random variable

$$f_Y(y) \delta y = P(Y \in [y, y + \delta y))$$



$$\frac{\delta y}{\delta x} = \frac{dg(x)}{dx} = g'(x) \leftarrow \text{derivative of } g(x) \text{ w.r.t } x$$

$$\delta y = g'(x) \delta x \quad \text{where } \delta x \text{ is small}$$

$$\delta y = g'(x) \delta x \quad \text{when } \delta x \text{ is small}$$

$$\int f(y) \delta y = P(Y \in [y, y + \delta y)) = P\left(X \in \left[\underbrace{g^{-1}(x)}_{x_g}, \underbrace{g^{-1}(x) + |g'(x)| \delta x}_{x_g + \delta x_g}\right)\right)$$

$$= P(X \in [x_g, x_g + \delta x_g))$$

$$\begin{aligned}
 &= P(X \in [x_g, x_g + |g'(x)| \delta x]) \\
 &= g'(x) P(X \in [x_g, x_g + \delta x]) \\
 &= g'(x) f_x(x_g) \delta x
 \end{aligned}$$

$\xleftrightarrow{g'(x)}$
 $|$
 $|$
 $x + \delta x$
 $\underbrace{\hspace{1cm}}_{g'(x) \delta x}$

$$\begin{aligned}
 &P(Y \in [y, y + \delta y]) \\
 &= P(\underbrace{g^{-1}(Y)}_X \in [g^{-1}(y), g^{-1}(y + \delta y)]) \\
 &= P(X \in [g^{-1}(y), g^{-1}(y + \delta y)])
 \end{aligned}$$

$Y = g(X)$
 $\approx X^2$
 $\frac{f_x(x)}{\text{given}}$
 $f_Y(y) = ?$

$g^{-1}(y + \delta y) \approx \text{Taylor Series expansion} \quad g^{-1}(y) + \underbrace{\frac{d}{dy} g^{-1}(y)}_{\text{}} \delta y$

+ ... higher order terms

$$\frac{d}{dy} g^{-1}(y) = \lim_{\delta y \rightarrow 0} \frac{g^{-1}(y + \delta y) - g^{-1}(y)}{\delta y}$$

$$\lim_{\delta y \rightarrow 0} \left(\frac{d}{dy} g^{-1}(y) \right) \delta y = g^{-1}(y + \delta y) - g^{-1}(y)$$

$$\lim_{\delta y \rightarrow 0} \quad \underline{g^{-1}(y + \delta y) = g^{-1}(y) + \left[\frac{d}{dy} g^{-1}(y) \right] \delta y}$$

first order
Taylor series
expansion

$$f_Y(y) \delta y$$

$$= P(Y \in [y, y + \delta y))$$

$$= P(X \in [g^{-1}(y), g^{-1}(y + \delta y)))$$

$$= P\left(X \in \left[g^{-1}(y), g^{-1}(y) + \underbrace{\frac{d}{dy} g^{-1}(y) \delta y}_{\delta x_y}\right]\right)$$

$$= P(X \in [g^{-1}(y), g^{-1}(y) + \delta x_y))$$

$$= f_X(g^{-1}(y)) \delta x_y$$

$$f_Y(y) \delta y = f_X(g^{-1}(y)) \left(\frac{d}{dy} g^{-1}(y) \right) \delta y$$

$$f_Y(y) = f_X(g^{-1}(y)) \left| \frac{d}{dy} g^{-1}(y) \right|$$

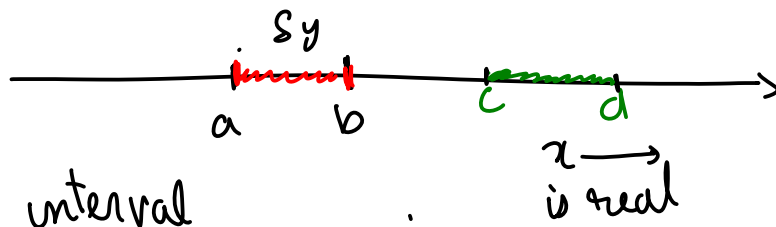
Transformation of
cont. RV

Example

$$y = g(x) = x^2 \Rightarrow g^{-1}(y) = +\sqrt{y}$$

$$\Rightarrow \frac{d}{dy} g^{-1}(y) = \frac{1}{2\sqrt{y}}$$

$$f_y(y) = f_x(\sqrt{y}) \left| \frac{1}{2\sqrt{y}} \right|$$



closed interval

$x \in [a, b]$ includes a, b

Open interval

$x \in (c, d)$ excludes c, d

Half closed, half open

$x \in [a, b)$
includes a
but excludes b

$x \in (a, b]$
includes b
but excludes a

Expectations and covariances

$$\mathbb{E}[f] = \sum_x p(x) f(x)$$

$$\mathbb{E}[f] = \int p(x) f(x) dx.$$

A RV X and a function of RV X , $f(X)$

Discrete RV

$$\mathbb{E}_X[f] = \sum_{x \in \Omega(X)} p(X=x) f(x)$$

Conditional expectation

Cont. RV

$$\mathbb{E}_X[f] = \int_{x \in \mathcal{R}(X)} \overset{\text{Prob density function}}{p(x)} f(x) dx$$

$$\mathbb{E}_x[f|y] = \sum_x p(x|y) f(x)$$

6-side dice
 $X = 1 \ 2 \ 3 \ 4 \ 5 \ 6$
 $p(X) = 1/6$
 unbiased die
 $M(X) = \begin{cases} \$5 & \text{if you get } X=2 \\ \$15 & \text{if you get } X=4 \\ 0 & \text{otherwise} \end{cases}$

$$\begin{aligned} \mathbb{E}_X[M(X)] &= \frac{5}{6} + \frac{15}{6} + 0 + 0 + 0 + 0 \\ &= \frac{20}{6} = \frac{10}{3} \end{aligned}$$

~~Expected value is area under the curve~~

Conditional expectation

$$\begin{aligned} \mathbb{E}_X[f(x) | Y=y] &= \sum_x p(X=x | Y=y) f(x) \\ &= \int_x p(x | Y=y) f(x) dx \end{aligned}$$

Law of iterated expectation

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

Sum rule

$$P(X) = \sum_y P(X, Y=y)$$

$$E_x[f(X)] = E_y[E_x[f(X) | Y=y]]$$

$$= \sum_y P(X | Y=y) P(Y=y)$$

$$\underbrace{E_x[f(X) | Y=y]}_{g(y)} = \underbrace{\sum_x f(x) P(X=x | Y=y)}_{g(y)}$$

Def of cond expectation

$$E_y[g(Y)] = \sum_y g(y) P(Y=y)$$

$$= \sum_y \left(\sum_x f(x) P(X=x | Y=y) \right) P(Y=y)$$

$$= \sum_y \left(\sum_x f(x) P(X=x, Y=y) \right)$$

$$= \sum_x f(x) \sum_y P(X=x, Y=y)$$

$$= \sum_x f(x) P(X=x)$$

$$= E_x[f(X)] = LHS$$

$$P(X=x | Y=y) = \frac{P(X=x, Y=y)}{P(Y=y)}$$

Variance

$$\mu = \mathbb{E}_x[f(x)]$$

$$\text{var}[f] = \mathbb{E} [(f(x) - \mathbb{E}[f(x)])^2]$$

$$\text{var}_x[f] = \mathbb{E}_x [(f(x) - \mu)^2]$$

Covariance

$$\begin{aligned} \text{cov}[x, y] &= \mathbb{E}_{x, y} [\{x - \mathbb{E}[x]\} \{y - \mathbb{E}[y]\}] \\ &= \mathbb{E}_{x, y} [xy] - \mathbb{E}[x]\mathbb{E}[y] \end{aligned}$$

$$\text{cov}[X, Y] = \mathbb{E}_{X, Y} [(X - \mu_X)(Y - \mu_Y)]$$

$$\mu_X = \mathbb{E}_X[X]$$

$$\mu_Y = \mathbb{E}_Y[Y]$$

Covariance in case of random vectors

$$\begin{aligned} \text{cov}[\mathbf{x}, \mathbf{y}] &= \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\{\mathbf{x} - \mathbb{E}[\mathbf{x}]\} \{\mathbf{y}^T - \mathbb{E}[\mathbf{y}^T]\}] \\ &= \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\mathbf{x}\mathbf{y}^T] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{y}^T]. \end{aligned}$$

$$\text{cov}(\underline{x}, \underline{y}) = \mathbb{E}_{\underline{x}, \underline{y}} [\underbrace{(\underline{x} - \underline{\mu}_x)(\underline{y} - \underline{\mu}_y)^T}_{\text{outer product}}]$$

$$\begin{aligned} \underline{x} &= \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \\ \underline{y} &= \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \end{aligned}$$

Inner product / Dot product

$$\underline{x}^T \underline{y} = \underline{x} \cdot \underline{y} = \underline{y}^T \cdot \underline{x} \in \mathbb{R}$$

Outer product

$$\underbrace{\underline{x}}_{n \times 1} \underbrace{\underline{y}^T}_{1 \times m} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} [y_1 \dots y_m]$$

$$\underline{\mu}_x = \mathbb{E}_{\underline{x}} [\underline{x}]$$

$$= \begin{bmatrix} \mathbb{E}_{x_1}[x_1] \\ \mathbb{E}_{x_2}[x_2] \\ \vdots \\ \mathbb{E}_{x_n}[x_n] \end{bmatrix}$$

$$= \begin{bmatrix} x_1 y_1 & x_1 y_2 & \dots & x_1 y_m \\ x_n y_1 & & & x_n y_m \end{bmatrix} \in \mathbb{R}^{n \times m}$$

$$\text{cov}(\underline{x}, \underline{y}) = E_{\underline{x}, \underline{y}} \left[\underbrace{(\underline{x} - \underline{\mu}_x)(\underline{y} - \underline{\mu}_y)^T}_{\text{outer product}} \right]$$

$$= \begin{bmatrix} \text{cov}(x_1, y_1) & \dots & \text{cov}(x_1, y_m) \\ \text{cov}(x_n, y_1) & \dots & \text{cov}(x_n, y_m) \end{bmatrix}$$

$$\begin{aligned} \text{cov}(\underline{x}) = \text{var}(\underline{x}) &= E_x \left[\overbrace{(\underline{x} - \underline{\mu}_x)}^{\mathbf{p}} \overbrace{(\underline{x} - \underline{\mu}_x)^T}^{\mathbf{p}^T} \right] \\ &= \begin{bmatrix} \text{var}(x_1) & \text{cov}(x_1, x_2) & \dots & \text{cov}(x_1, x_n) \\ \vdots & \text{var}(x_2) & & \\ \vdots & & \ddots & \\ \text{cov}(x_n, x_1) & & & \text{var}(x_n) \end{bmatrix} \end{aligned}$$

var(x) is always positive ^{semi} definite

$$\underline{z}^T \text{var}(\underline{x}) \underline{z} \geq 0 \quad \text{for all } \underline{z}$$

$$\underbrace{\underline{z}^T \begin{bmatrix} \mathbf{p} & \mathbf{p}^T \\ \mathbf{p} & \mathbf{p}^T \end{bmatrix} \underline{z}}_{\text{outer product}} = (\mathbf{p}^T \underline{z})^2 > 0$$

Bayes theorem : Posterior, likelihood and prior

$$p(\underline{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\underline{w})p(\underline{w})}{p(\mathcal{D})}$$

weights
Dataset

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

$$\text{posterior} \propto \text{likelihood} \times \text{prior}$$

$$\text{Model} = \hat{y} = f(\underline{x}; \underline{w})$$

prediction
Hypothesis
hidden
input feature
Parameters or weight } Unknowns

$$p(\mathcal{D}|\underline{w}) = \underline{\text{likelihood of observation}}$$

observed

$$p(\underline{w}) = \text{A prior on } \underline{\text{the thing to be estimated}}$$

$$p(\mathcal{D}) = \text{Evidence}$$

$$p(\underline{w}|\mathcal{D}) = \underline{\text{Posterior}}$$

$$P(\mathcal{D}, \underline{w}) = P(\underline{w}, \mathcal{D})$$

$$P(\mathcal{D}|\underline{w})p(\underline{w}) = P(\underline{w}|\mathcal{D})p(\mathcal{D})$$

$$\frac{p(\mathcal{D}|\underline{w})p(\underline{w})}{p(\mathcal{D})} = \underline{p(\underline{w}|\mathcal{D})}$$

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

Metrics of binary classification

		Predicted condition	
Total population = P + N		Positive (PP)	Negative (PN)
Actual condition	Positive (P)	Recall = TPR True positive (TP) $TPR = TP/P$	False negative (FN) Type II error
	Negative (N)	False positive (FP) Type I error	True negative (TN)

Precision = TP/PP

Accuracy = $(TP+TN)/(P+N)$

- ① Precision : TP / PP
- ② Recall : TP / P
- ③ Accuracy : $(TP+TN)/(P+N)$
- ④ F1-score : Harmonic mean of Precision and Recall

Detected by test

	cancer	normal
True cancer	0	1000
True normal	1	0

Confusion Matrix

$loss(1, 1)$ $loss(0, 1)$
 $loss(1, 0)$ $loss(0, 0)$

$\hat{y} = f(x) = \hat{y}(x)$

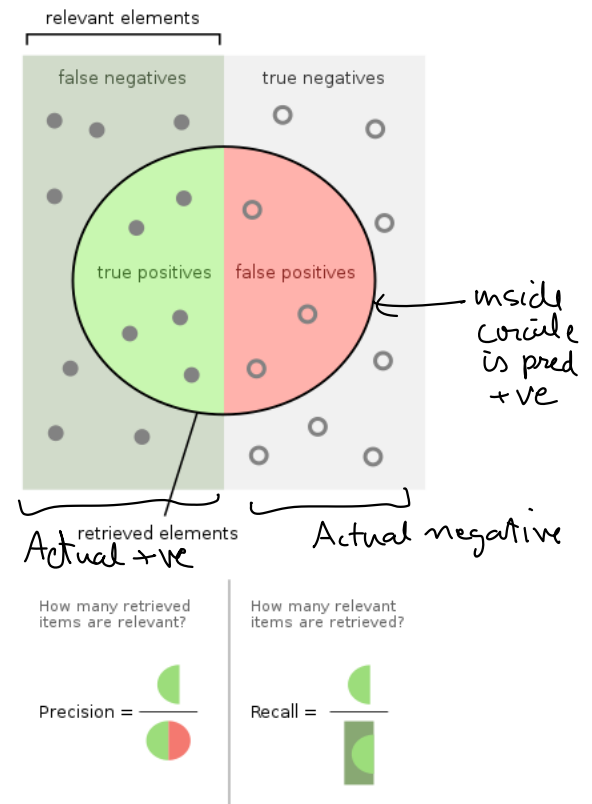
Recall

- True Positive Rate:** $TPR = \mathbb{P}[\hat{Y}(X) = 1 | Y = 1]$. Also known as power, sensitivity, probability of detection, or recall.
- False Negative Rate:** $FNR = 1 - TPR$. Also known as type II error or probability of missed detection.
- False Positive Rate:** $FPR = \mathbb{P}[\hat{Y}(X) = 1 | Y = 0]$. Also known as size or type I error or probability of false alarm.
- True Negative Rate** $TNR = 1 - FPR$, the probability of declaring $\hat{Y} = 0$ given $Y = 0$. This is also known as specificity.

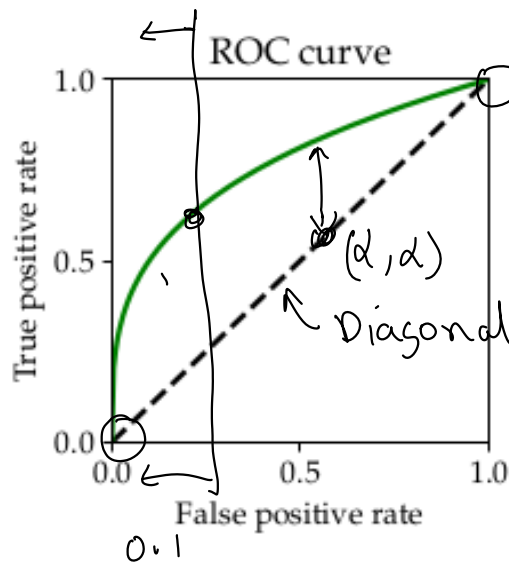
$$TPR = \frac{TP}{P} = \frac{P(\hat{y}=1, y=1)}{P(y=1)} = P(\hat{y}=1 | y=1)$$

F1-score is the harmonic mean of Precision and Recall

F1-score = $2 * \text{Recall} * \text{Precision} / (\text{Precision} + \text{Recall})$



Receiver Operating characteristics



$$P(\hat{Y}(x)=1) = \alpha$$

$$TPR = P(\hat{Y}(x)=1 \mid Y=1) = \alpha$$

$$FPR = P(\hat{Y}(x)=1 \mid Y=0) = \alpha$$

Diagonal of the ROC curve is a trivial Randomized classifier

Definition 1. We define the risk associated with \hat{Y} to be

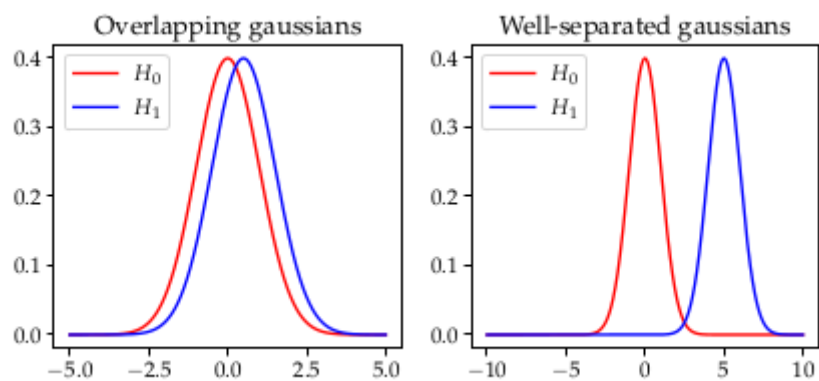
$$R[\hat{Y}] := \mathbb{E}[\text{loss}(\hat{Y}(X), Y)] .$$

Here, the expectation is taken jointly over X and Y .

Lemma 1. We claim that the optimal predictor is given by

$$\hat{Y}(x) = \mathbb{1} \left\{ \mathbb{P}[Y = 1 \mid X = x] \geq \frac{\text{loss}(1, 0) - \text{loss}(0, 0)}{\text{loss}(0, 1) - \text{loss}(1, 1)} \mathbb{P}[Y = 0 \mid X = x] \right\}$$

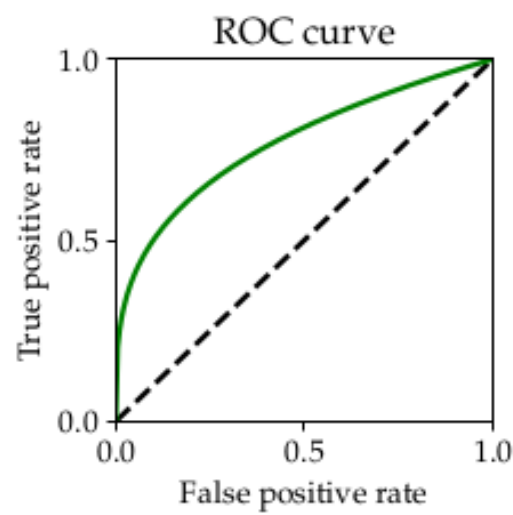
Likelihood ratio and likelihood ratio test



Maximum a posteriori predictor

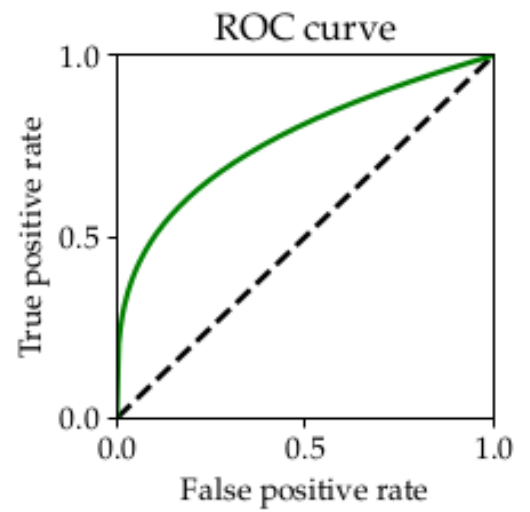
Maximum likelihood predictor

Receiver operating characteristic (ROC curve)



Neyman-Pearson Lemma

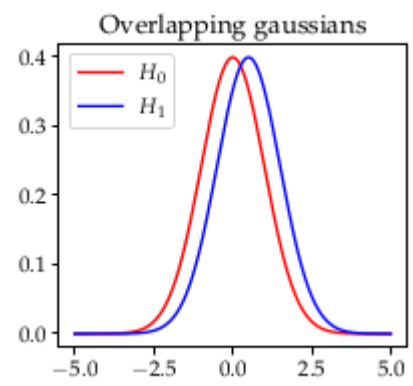
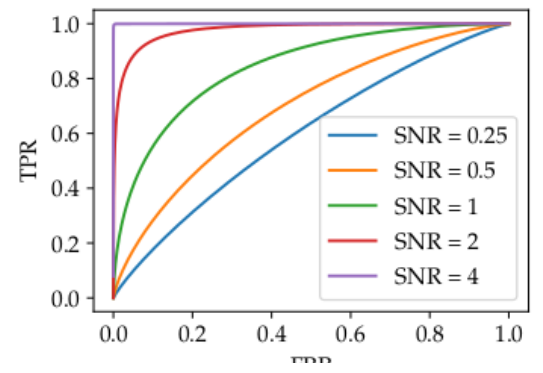
Lemma 2. Neyman-Pearson Lemma. Suppose the likelihood functions $p(x|y)$ are continuous. Then the optimal probabilistic predictor that maximizes TPR with an upper bound on FPR is a deterministic likelihood ratio test



Proposition 1. The points $(0, 0)$ and $(1, 1)$ are on the ROC curve.

Proposition 2. The ROC must lie above the main diagonal.

Proposition 3. The ROC curve is concave.



Decision Theory

Consider, for example, a medical diagnosis problem in which we have taken an X-ray image of a patient, and we wish to determine whether the patient has cancer or not. In this case, the input vector \mathbf{x} is the set of pixel intensities in the image, and output variable t will represent the presence of cancer, which we denote by the class C_1 , or the absence of cancer, which we denote by the class C_2 . We might, for instance, choose t to be a binary variable such that $t = 0$ corresponds to class C_1 and $t = 1$ corresponds to class C_2 . The general inference problem then involves determining the joint distribution $p(\mathbf{x}, C_k)$, or equivalently $p(\mathbf{x}, t)$, which gives us the most complete probabilistic description of the situation.

We would like this choice to be optimal in some appropriate sense. This is the decision step, and it is the subject of decision theory to tell us how to make optimal decisions given the appropriate probabilities.

$$p(C_k|\mathbf{x}) = \frac{p(\mathbf{x}|C_k)p(C_k)}{p(\mathbf{x})}.$$

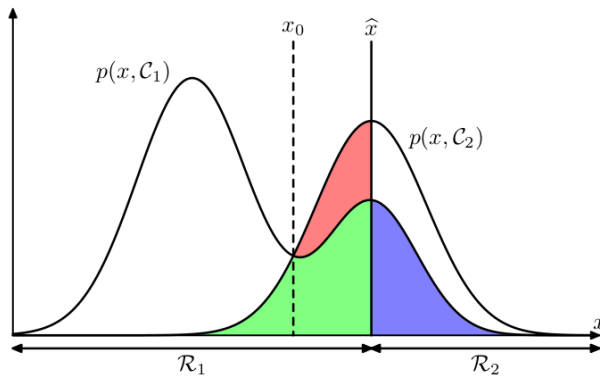


Figure 1.24 Schematic illustration of the joint probabilities $p(x, C_k)$ for each of two classes plotted against x , together with the decision boundary $x = \hat{x}$. Values of $x \geq \hat{x}$ are classified as class C_2 and hence belong to decision region \mathcal{R}_2 , whereas points $x < \hat{x}$ are classified as C_1 and belong to \mathcal{R}_1 . Errors arise from the blue, green, and red regions, so that for $x < \hat{x}$ the errors are due to points from class C_2 being misclassified as C_1 (represented by the sum of the red and green regions), and conversely for points in the region $x \geq \hat{x}$ the errors are due to points from class C_1 being misclassified as C_2 (represented by the blue region). As we vary the location \hat{x} of the decision boundary, the combined areas of the blue and green regions remains constant, whereas the size of the red region varies. The optimal choice for \hat{x} is where the curves for $p(x, C_1)$ and $p(x, C_2)$ cross, corresponding to $\hat{x} = x_0$, because in this case the red region disappears. This is equivalent to the minimum misclassification rate decision rule, which assigns each value of x to the class having the higher posterior probability $p(C_k|x)$.

$$\begin{aligned} p(\text{mistake}) &= p(\mathbf{x} \in \mathcal{R}_1, C_2) + p(\mathbf{x} \in \mathcal{R}_2, C_1) \\ &= \int_{\mathcal{R}_1} p(\mathbf{x}, C_2) d\mathbf{x} + \int_{\mathcal{R}_2} p(\mathbf{x}, C_1) d\mathbf{x}. \end{aligned}$$

$$\begin{aligned}
 p(\text{correct}) &= \sum_{k=1}^K p(\mathbf{x} \in \mathcal{R}_k, \mathcal{C}_k) \\
 &= \sum_{k=1}^K \int_{\mathcal{R}_k} p(\mathbf{x}, \mathcal{C}_k) \, \mathrm{d}\mathbf{x}
 \end{aligned}$$

ROC curves