# Machine Learning Algorithms

**List of concepts**
1   Dot product in Feature space
2   State space and Feature space
3   Correlation and Convolution
4   Markovian vs Martingale

**List of learning algorithms**

| | | | |
|---|---|---|---|
| 1.1 Bayes rule | | | |
| 1.2 Naïve Bayes classifier | | | |
| 1.3 Maximum likelihood | | *X is data* | *μΣ are parameters* |
| 1.4 Maximum aposteria | | *X is data* | *μΣ are parameters* |
| 1.5 Linear regression *(AX-B in ASM)* | *max likelihood* | *A,B are data* | *X is parameter* |
| 2.1 Logistic regression | *max likelihood* | *X,Y are data* | *θ is parameter* |
| 2.2 Expectation maximization | *max likelihood* | *X(Z) are data* | *θ is parameter* |
| 2.3 Gaussian mixture model | *max likelihood* | *X(Z) are data* | *μΣ are parameters* |
| 2.4 K mean clustering | *min distance* | *X(Z) are data* | *μ is parameter* |
| 2.5 K nearest neightbour | *min distance* | | |
| 2.6 KL divergence | *min divergence* | *X is data* | *θ is parameter* |
| 2.7 Decision tree | *min entropy* | *X,Y are data* | |
| 3.1 Kalman filter | *max likelihood* | *A,B are data* | *X is parameter* |
| 3.2 Principal component analysis | *variance analysis* | *X is data* | |
| 3.3 Linear discriminant analysis | *variance analysis* | *X,Y are data* | |
| 3.4 Neural network (RNN, CNN) | *min empirical risk* | *X,Y are data* | *w is parameter* |
| 3.5 Support vector machine | *min structural risk* | *X,Y are data* | *α is parameter* |
| 3.6 Generative adversarial network | *game theory (Oligopoly)* | | |
| 4.1 Reinforcement learning | *Markov chain* | | |
| 4.2 Hidden Markov model | *Markov chain* | | |

---

Generative vs Discriminative classifers

Generative classifers models distribution of each class without considering interclass differences, and find the nearest in inspection. Discriminative classifers model the differences or boundary between classes without considering distribution of each class. There is no *decision boundary* between classes for generative classifiers, there is no *uniclass distribution model* for discriminative classifier.

| *Generative classifer models P(X,y)* | *Discriminative classifier models P(X|y)* |
|---|---|
| *class 1*     *P(X, y=1)* | |
| *class 2*     *P(X, y=2)* | |
| *…* | |
| *example : Naïve Bayes, GMM, GAN* | *example : logestic regression, NN, SVM* |

**Summary** `9*3,3*3,3*10,3*4 + 354113333`

| 9*3 | | objective | data | para | content |
|---|---|---|---|---|---|
| 2.1 | LR | *max likelihood* | *X,Y* | θ | 3 = *define L(θ), max L(θ), linear case* |
| 2.2 | EM | *max likelihood* | *X(Z)* | θ | 5 = *define convex, log, define L'(θ), max L'(θ), algo* |
| 2.3 | GMM | *max likelihood* | *X(Z)* | $\mu\Sigma$ | 4 = *define pdf, Estep, Mstep, weight matrix* |
| 2.4 | Kmean | *min distance* | *X(Z)* | $\mu$ | 1 |
| 2.5 | KNN | *min distance* | | | 1 |
| 2.6 | KLdiv | *min divergence* | *X* | θ | 3 = *define entropy, define divergence, max likelihood* |
| 2.7 | Decision tree | *min entropy* | *X,Y* | | 3 = *inspection, construct tree, train tree* |
| 3.4 | NN | *min empirical risk* | *X,Y* | *w* | 3 = *define topology, algorithm, proof* |
| 3.5 | SVM | *min structural risk* | *X,Y* | $\alpha$ | 3 = 10,10,10 |

Comparison among Gauss Newton / logistic regression and neural network

| 3*3 | objective | matrix | algorithm |
|---|---|---|---|
| *Gauss Newton* | $\min\limits_{\Delta X}(F(A\mid X^{(t)})+J_F\Delta X-B)^T\Sigma^{-1}(F(A\mid X^{(t)})+J_F\Delta X-B)$ | $F(A\mid X),J_F$ | *Newton Raphson* |
| *logistic regression* | $\max\limits_{\theta}\ln\prod_{n=1}^{N}[s(f(x_n\mid\theta))]^{y_n}[1-s(f(x_n\mid\theta))]^{1-y_n}$ | $f(x_n\mid\theta),J_F$ | *grad ascent* $\Delta\theta=J_F^T E$ |
| *neural network* | $\min\limits_{W}(F(X\mid W)-Y)^T(F(X\mid W)-Y)$ | $F(X\mid w)\ only$ | *grad descent* $\Delta w=-s\partial_w E$ |

Definition of *F* and *J* in Gauss Newton, please modify accordingly for logistic regression and neural network :

$$F(A\mid X)=\begin{bmatrix}f(A_1\mid X)\\f(A_2\mid X)\\f(A_3\mid X)\\...\\f(A_N\mid X)\end{bmatrix}\qquad J_F=\begin{bmatrix}\partial f(A_1\mid X)/\partial x_1 & \partial f(A_1\mid X)/\partial x_2 & \partial f(A_1\mid X)/\partial x_3 & ... & \partial\partial f(A_1\mid X)/\partial x_M\\\partial f(A_2\mid X)/\partial x_1 & \partial f(A_2\mid X)/\partial x_2 & \partial f(A_2\mid X)/\partial x_3 & ... & \partial f(A_2\mid X)/\partial x_M\\\partial f(A_3\mid X)/\partial x_1 & \partial f(A_3\mid X)/\partial x_2 & \partial f(A_3\mid X)/\partial x_3 & ... & \partial f(A_3\mid X)/\partial x_M\\... & ... & ... & ... & ...\\\partial f(A_N\mid X)/\partial x_1 & \partial f(A_N\mid X)/\partial x_2 & \partial f(A_N\mid X)/\partial x_3 & ... & \partial f(A_N\mid X)/\partial x_M\end{bmatrix}$$

| 3*10 | | SVC | | SVC | | SVR |
|---|---|---|---|---|---|---|
| *#arguments in step1 obj-fct* | 2 | $w,b$ | 3 | $w,b,\xi$ | 4 | $w,b,\xi^+,\xi^-$ |
| *#constraints in step1* | 1 | $y_n f(x_n)\geq 1$ | 2 | $y_n f(x_n)\geq ...$ , $\xi_n\geq 0$ | 4 | $f(x_n)\pm\varepsilon...\geq y_n\pm\xi_n^{\pm}$ , $\xi_n^{\pm}\geq 0$ |
| *#arguments in step2 obj-fct* | 3=21 | $w,b,\alpha$ | 5=32 | $w,b,\xi,\alpha,\beta$ | 8=44 | $w,b,\xi^{\pm},\alpha^{\pm},\beta^{\pm}$ |
| *#constraints in step2* | 1 | $\alpha\geq 0$ | 2 | $\alpha\geq 0$ , $\beta\geq 0$ | 4 | $\alpha^{\pm}\geq 0$ , $\beta^{\pm}\geq 0$ |
| *#arguments in step3 obj-fct* | 1 | $\alpha$ | 1 | $\alpha$ | 1 | $A$ |
| *#constraints in step3* | 2 | $\alpha\geq 0$ , $y^T\alpha=0$ | 2 | $\alpha\in[0,C]$ , $y^T\alpha=0$ | 2 | $A\in\pm C$ , $l^T A=0$ |
| *y in Lag 2nd order term* | | *yes* | | *yes* | | - |
| *y in Lag 1st order term?* | | - | | - | | *yes* |
| *y in linear constraint?* | | *yes* | | *yes* | | - |
| *y in w vector?* | | *yes* | | *yes* | | - |

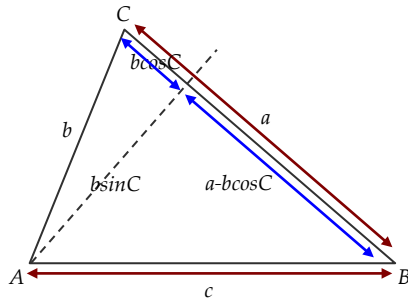| 3*4 | | | | | |
|---|---|---|---|---|---|
| *SVC* | $\max\limits_{\alpha}-\dfrac{1}{2}\alpha^T H\alpha+l^T\alpha$ | $st\ \alpha\geq 0\ and\ y^T\alpha=0$ | $where\ H_{nm}=y_n y_m x_n x_m^T$ | $and\ \alpha=\begin{bmatrix}\alpha_1\\\alpha_2\\\alpha_3\\...\end{bmatrix}$ | |
| *SVC* | $\max\limits_{\alpha}-\dfrac{1}{2}\alpha^T H\alpha+l^T\alpha$ | $st\ 0\leq\alpha\leq C\ and\ y^T\alpha=0$ | $where\ H_{nm}=y_n y_m x_n x_m^T$ | $and\ \alpha=\begin{bmatrix}\alpha_1\\\alpha_2\\\alpha_3\\...\end{bmatrix}$ | |
| *SVR* | $\max\limits_{\alpha^{\pm}}-\dfrac{1}{2}\alpha^T H\alpha+y^T\alpha-l^T\alpha*\varepsilon$ | $st\ 0\leq\alpha\leq C\ and\ y^T\alpha=0$ | $where\ H_{nm}=x_n x_m^T$ | $and\ \alpha=\begin{bmatrix}\alpha_1^+-\alpha_1^-\\\alpha_2^+-\alpha_2^-\\\alpha_3^+-\alpha_3^-\\...\end{bmatrix}$ | $\alpha*=\begin{bmatrix}\alpha_1^++\alpha_1^-\\\alpha_2^++\alpha_2^-\\\alpha_3^++\alpha_3^-\\...\end{bmatrix}$ |

## 1 Dot product in Feature space

Dot product of 2 feature vectors measures their cosine similarity. Let's prove it starting from Pythagoras theorem and law of cosine.

- *points $ACB_0$ are collinear and*
- *points $BCA_0$ are collinear only if C is a right angle*
- *$ABA_1$ has same area as $ACA_2$ (same sides and included amgle)*
- *$ACA_1$ has same area as $AC_1A_2$ (shifting B to C, shifting C to $C_1$)*
- *$ACA_0A_1$ has same area as $AC_1C_2A_2$ (doubling triangles into square & rectangle)*
- *$BCB_0B_1$ has same area as $BC_1C_2B_2$ (repeating all above steps for RHS)*
- *$AB^2 = AC^2 + BC^2$ (sum of the two results above)*

or     $c^2 = b^2 + a^2$

After that, we will generalise Pythagoras theorem to non-right angle triangle.

$$
\begin{aligned}
c^2 &= (b\sin C)^2 + (a - b\cos C)^2 \\
&= b^2\sin^2 C + a^2 - 2ab\cos C + b^2\cos^2 C \\
&= b^2 + a^2 \underbrace{- 2ab\cos C}_{\substack{extra-team \\ comparing \\ Pythagoras}}
\end{aligned}
$$

*The extra term vanishes if C is right angle.*

This is called the law of cosine. Now we can show why dot product is equivalent to projection, and thus related to cosine similarity. Given vector $A, B \in \mathfrak{R}^M$, and $C = A - B$, then we have :

$$
\begin{aligned}
c^2 &= \vec{C} \cdot \vec{C} \\
&= (\vec{A} - \vec{B}) \cdot (\vec{A} - \vec{B}) \\
&= \vec{A} \cdot \vec{A} - 2\vec{A}\vec{B} + \vec{B} \cdot \vec{B} \\
&= b^2 + a^2 - 2\vec{A}\vec{B} \\
&= b^2 + a^2 - 2ab\cos C \qquad \text{using law of cosine}
\end{aligned}
$$

| Recall : | | |
|---|---|---|
| $\vec{A} \cdot \vec{B}$ | = | *dot product* |
| $\vec{A} \cdot \vec{B}/\lvert\vec{B}\rvert$ | = | *projection of A on B* |
| $\vec{A} \cdot \vec{B}/(\lvert\vec{A}\rvert\lvert\vec{B}\rvert)$ | = | *cosine similarity* |

Therefore we have :

$$
\vec{A}\vec{B} = ab\cos C = \lvert\vec{A}\rvert\lvert\vec{B}\rvert\cos(\angle AOB) \qquad \text{where O is the origin}
$$

## 2 State space and Feature space

Discrete state space is represented as a set, $S = \{S_0, S_1, \ldots, S_N\}$, while continuous state space is usually represented as multidimension feature space, $F \in \mathfrak{R}^M$ under some boundary constraints. Here is the correspondence :

## 3 Correlation and Convolution

Correlation is defined as inner product :

$$
\begin{aligned}
c(t) &= \int f(s)g(s+t)ds \\
c(0) &= \int f(s)g(s)ds \qquad\qquad \textit{i.e. dot product of functions}
\end{aligned}
$$

Convolution is defined as :

$$
c(t) = \int f(s)g(t-s)ds
$$

## 4 Markovian vs Martingale

Markovian system is a memoryless system, probability of future states conditional on current information depends on current state only, while historical states are not involved, in other words :

$$
\begin{aligned}
&\Pr(X_{n+1}=x\,|\,X_n,...,X_1,X_0) = \Pr(X_{n+1}=x\,|\,X_n) &&\textit{for discrete space} \\
\textit{or}\quad &\Pr(X_T \in A\,|\,I_t) = \Pr(X_T \in A\,|\,X_t) &&\textit{for continuous space}
\end{aligned}
$$

Examples of Markovian process include : Markov chain which is usually represented as a directed graph, chess and Black Scholes :

$$
dX_t = \mu X_t dt + \sigma X_t dz_t \qquad\qquad \textit{depending on current } X_t \textit{ only}
$$

Non Markovian process include autoregressive time series which can be generalized into *ARMA* and *GARCH*. *AR* model is :

$$
X_{n+1} = \sum_{m=0}^{M} w_m X_{n-m} + \varepsilon_{n+1} \qquad\qquad \textit{depending on historical states, declare M extra variables for history}
$$

Given *M+1* observation set, the parameter vector $w$ can be solved by least square, resulting in the Yule Walker equation.

Martingale process is defined as a process, which has future expectation conditional on current information equals to current value.

$$
\begin{aligned}
&E(X_{n+1}\,|\,X_n,...,X_1,X_0) = X_n &&\textit{for discrete space} \\
\textit{or}\quad &E(X_T\,|\,I_t) = X_t &&\textit{for continuous space}
\end{aligned}
$$

Example of martingale process include risk neutral Black Scholes with risk free discounting. Martingale concept is important in the fundalmental theorem of asset pricing, from which risk neutral pricing is derived. Continuous martingale is equivalent to :

$$
\begin{aligned}
&E(X_{t+dt}\,|\,I_t) = X_t \\
\Rightarrow\quad &E(X_t + dX_t\,|\,I_t) = X_t \\
\Rightarrow\quad &E(dX_t\,|\,I_t) = X_t - E(X_t\,|\,I_t) \\
&\qquad\qquad\quad = 0 \qquad\qquad \textit{hence it implies process with no drift}
\end{aligned}
$$

Comparison

Therefore Markovian is defined using conditional probability, whereas martingale is defined using conditional expectation, besides, the former has *RHS* depending on current state, while the latter has *RHS* equivalent to current state. Comparison in discrete space :

$$
\begin{aligned}
&\Pr(X_{n+1}=x\,|\,X_n,...,X_1,X_0) = \Pr(X_{n+1}=x\,|\,X_n) &&\textit{for Markovian} \\
&E(X_{n+1}\,|\,X_n,...,X_1,X_0) = X_n &&\textit{for Martingale}
\end{aligned}
$$

There are *4* combinations :

| | Martingale | non martingale |
|---|---|---|
| *Markovian* | $dX_t = \sigma f(X_t)dz_t$ | $dX_t = \mu dt + \sigma f(X_t)dz_t$ |
| *non Markovian* | $dX_t = \sigma f(X_s)dz_t \quad \forall s < t$ | $dX_t = \mu dt + \sigma f(X_s)dz_t \quad \forall s < t$ |

## 1.1 Bayes rule

Bayes rule relates posterior probability with prior probability.

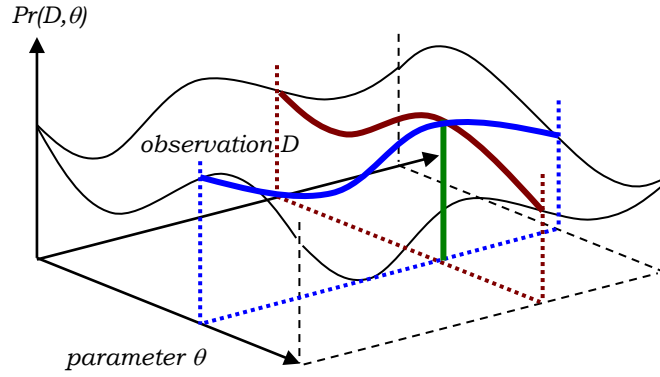$$\Pr(\theta \mid D) \quad = \quad \frac{\Pr(D \mid \theta)\,\Pr(\theta)}{\Pr(D)}$$

where
| | | | | |
|---|---|---|---|---|
| | $\theta$ | = | hypothesis | (i.e. model parameters) |
| | $D$ | = | evidence | (i.e. noisy observations) |
| | $\Pr(\theta)$ | = | prior probability | (i.e. probability of hypothesis without observation) |
| | $\Pr(\theta \mid D)$ | = | posterior probability | (i.e. probability of hypothesis given some observations) |
| | $\Pr(D \mid \theta)$ | = | likelihood | (i.e. probability of having observations given the hypothesis) |

In terms of joint distribution, we have :

$$\Pr(\theta \mid D) \quad = \quad \frac{\Pr(D,\theta)}{\Pr(D)}$$

$$\Pr(D \mid \theta) \quad = \quad \frac{\Pr(D,\theta)}{\Pr(\theta)}$$

where
| | | | |
|---|---|---|---|
| | $\Pr(\theta)$ | = | volume under blue curve |
| | $\Pr(D)$ | = | volume under red curve |
| | $\Pr(D,\theta)$ | = | volume of green pillar |



*Pr(D,θ)*

*observation D*

*parameter θ*

Two methods of model parameter estimation :

Maximum likelihood $\qquad \underset{\theta}{\arg\max}\ln\Pr(D \mid \theta)$

Maximum aposterior $\qquad \underset{\theta}{\arg\max}\ln\Pr(\theta \mid D) \quad = \quad \underset{\theta}{\arg\max}\ln\dfrac{\Pr(D \mid \theta)\,\Pr(\theta)}{\Pr(D)}$

$$= \quad \underset{\theta}{\arg\max}\ln[\Pr(D \mid \theta)\,\Pr(\theta)] \qquad \textit{where Pr is pmf or pdf}$$

## 1.2 Naïve Bayes Classifier

Naïve Bayes is a maximum-joint-probability-classifier, which breaks down the joint-probability into a chain of single-dimensional conditional-probabilities. Given $N$ training vectors having $M$ attributes together with class label $x_n \in \Re^M$ $y_n \in [1,K]$ $\forall n \in [1,N]$, there are $K$ possible classes, it classifies test vector $x \in \Re^M$ naively by considering conditional-probabilities in each dimension separately :

$$y \quad = \quad \underset{k}{\arg\max}\Pr(x, y = k)$$

$$= \quad \underset{k}{\arg\max}\Pr(x_{[1]} \mid x_{[2:M]}, y = k) \times \Pr(x_{[2:M]}, y = k) \qquad \textit{index [m] denotes attribute m}$$

$$= \quad \underset{k}{\arg\max}\Pr(x_{[1]} \mid x_{[2:M]}, y = k) \times \Pr(x_{[2]} \mid x_{[3:M]}, y = k) \times \Pr(x_{[3:M]}, y = k)$$

$$= \quad \underset{k}{\arg\max}\Pr(x_{[1]} \mid x_{[2:M]}, y = k) \times \Pr(x_{[2]} \mid x_{[3:M]}, y = k) \times .... \times r(x_{[M-1:M]} \mid x_{[M]}, y = k) \times \Pr(x_{[M]} \mid y = k) \times \Pr(y = k)$$

$$= \quad \underset{k}{\arg\max}\prod_{m=1}^{M}\underbrace{\Pr(x_{[m]} \mid y = k)}_{cond-prob-1D-attr}\;\Pr(y = k) \qquad \textit{if all M attributes are independent}$$

*For continuous case*
We have to make model-assumptions *(such as 1D Gaussian)* on conditional-probability for each *1D* attribute, we estimate the model parameters using training dataset (and some algorithms). Finally we can classify incoming test cases using the above formula.

*For discrete case*
If all the attributes are discrete values, then we can find the conditional probabilities from training set simply by counting.

$$\Pr(x_{[m]} \mid y = k) \quad = \quad \frac{\sum_{n=1}^{N}1(x_{n[m]} = v, y_n = k)}{\sum_{n=1}^{N}1(y_n = k)} \qquad \textit{for all } v \in [\textit{possible values of attribute m}]$$

$$\textit{for all } m \in [1,M] \textit{ and for all } k \in [1,K]$$

Thus the learnt record is a *3D* tensor.

## 1.3 Maximum Likelihood

Lets try maximum likelihood estimation with multivariate Gaussian. Given a set of *i.i.d.* (independent-identically distributed) data points $X = \{x_1, x_2, x_3, \ldots, x_N\}$ in $\Re^M$, drawn from $\mathcal{N}(x; \mu, \Sigma)$, suppose all data $x_n$ and $\mu$ are row vectors, then log likelihood function is :

$$\ln P(X \mid \mu, \Sigma) = \ln \left[ \prod_{n=1}^{N} \frac{1}{\sqrt{(2\pi)^{rank} \det(\Sigma)}} \exp\left( -\frac{(x_n - \mu)\Sigma^{-1}(x_n - \mu)^T}{2} \right) \right]$$

*we then apply the "trace-trick"*

$$= -\frac{N}{2}\ln\det(\Sigma) - \frac{1}{2}\sum_{n=1}^{N}(x_n - \mu)\Sigma^{-1}(x_n - \mu)^T + const \qquad (1)$$

$$= -\frac{N}{2}\ln\det(\Sigma) - \frac{1}{2}\sum_{n=1}^{N}tr((x_n - \mu)\Sigma^{-1}(x_n - \mu)^T) + const$$

*since trace(scalar) = scalar*

$$= -\frac{N}{2}\ln\det(\Sigma) - \frac{1}{2}\sum_{n=1}^{N}tr(\Sigma^{-1}(x_n - \mu)^T(x_n - \mu)) + const$$

*since trace is invariant to cyclic permutation*

$$= -\frac{N}{2}\ln\det(\Sigma) - \frac{1}{2}tr(\Sigma^{-1}\sum_{n=1}^{N}(x_n - \mu)^T(x_n - \mu)) + const \qquad (2)$$

Differentiate log likelihood in *(1)* wrt $\mu$ and setting it to zero, we have :

$$0 = -\frac{1}{2}\sum_{n=1}^{N}(-2)\Sigma^{-1}(x_n - \mu)^T$$

*since $\partial_X(AX - B)\Sigma(AX - B)^T = 2A\Sigma(AX - B)^T$*

$$\mu = \frac{1}{N}\sum_{n=1}^{N}x_n$$

Differentiate log likelihood in *(2)* wrt $\Sigma^{-1}$ and setting it to zero, we have :

$$\partial_{\Sigma^{-1}}\left(\frac{N}{2}\ln\det(\Sigma)\right) = \partial_{\Sigma^{-1}}\left(-\frac{1}{2}tr(\Sigma^{-1}\sum_{n=1}^{N}(x_n - \mu)^T(x_n - \mu))\right)$$

$$= -\frac{1}{2}\sum_{n=1}^{N}((x_n - \mu)^T(x_n - \mu))^T$$

*since $\partial_A Tr(AB) = \partial_A Tr(BA) = B^T$*

$$= -\frac{1}{2}\sum_{n=1}^{N}(x_n - \mu)^T(x_n - \mu)$$

*since $(AB)^T = B^T A^T$*

$$\partial_\Sigma\left(\ln\det(\Sigma)\right)\partial_{\Sigma^{-1}}\Sigma = -\frac{1}{N}\sum_{n=1}^{N}(x_n - \mu)^T(x_n - \mu)$$

$$\Sigma^{-T}\partial_{\Sigma^{-1}}\Sigma = -\frac{1}{N}\sum_{n=1}^{N}(x_n - \mu)^T(x_n - \mu)$$

*since $\partial_A\left(\ln\det(A)\right) = (A^{-1})^T = A^{-T}$*

$$\Sigma^{-T}(-\Sigma\Sigma) = -\frac{1}{N}\sum_{n=1}^{N}(x_n - \mu)^T(x_n - \mu)$$

*since $\partial_A A^{-1} = -A^{-1}A^{-1}$ and $\partial_{A^{-1}}A = -AA$*

$$\Sigma = \frac{1}{N}\sum_{n=1}^{N}(x_n - \mu)^T(x_n - \mu)$$

*result is intuitive*

*Remark for matrix differentiation*

*Lets prove $\partial_A\left(\ln\det(A)\right) = (A^{-1})^T = A^{-T}$ . Consider derivative wrt element $a_{ij}$ :*

Recall for matrix X with row data
$X^T X$ is covariance matrix
$X X^T$ is dot-product matrix

$$\partial_{a_{ij}}\left(\ln\det(A)\right) = (\det(A))^{-1}\partial_{a_{ij}}\det(A)$$

$$= (\det(A))^{-1}cofactor(A)_{ij}$$

$$= (A^{-1})_{ji} \qquad \Rightarrow \quad \partial_A\left(\ln\det(A)\right) = (A^{-1})^T = A^{-T}$$

*Lets prove $\partial_A A^{-1} = -A^{-1}A^{-1}$ .*

$$I = AA^{-1}$$

$$\partial_A I = A(\partial_A A^{-1}) + (\partial_A A)A^{-1}$$

$$0 = A(\partial_A A^{-1}) + A^{-1} \qquad \Rightarrow \quad \partial_A A^{-1} = -A^{-1}A^{-1}$$

## 1.4 Maximum Aposterior

If we are given prior knowledge about $\mu$ and $\Sigma$, we can incorporate them into maximum aposterior.

$$\ln[P(X\mid\mu,\Sigma)\Pr(\mu,\Sigma)] \quad = \quad \ln\left[\prod_{n=1}^{N}\frac{1}{\sqrt{(2\pi)^{rank}\det(\Sigma)}}\exp\left(-\frac{(x_n-\mu)\Sigma^{-1}(x_n-\mu)^T}{2}\right)\right]+\ln P_{prior}(\mu,\Sigma)$$

Suppose we know in prior that, parameter $\mu$ follows a $M$ dimensional Gaussian centred at $\mu_0$, having covariance matrix $\Sigma_0$, whereas nothing is known about $\Sigma$, then we have (recall that both $\mu$ and $\mu_0$ are row vectors) :

$$\ln[P(X\mid\mu,\Sigma)\Pr(\mu,\Sigma)] \quad = \quad \ln\left[\prod_{n=1}^{N}\frac{1}{\sqrt{(2\pi)^{rank}\det(\Sigma)}}\exp(-\frac{1}{2}(x_n-\mu)\Sigma^{-1}(x_n-\mu)^T)\right]+\ln\left[\frac{1}{\sqrt{(2\pi)^M\det(\Pi)}}\exp(-\frac{1}{2}(\mu-\mu_0)\Sigma_0^{-1}(\mu-\mu_0)^T)\right]$$

$$= \quad -\frac{N}{2}\ln\det(\Sigma)-\frac{1}{2}\sum_{n=1}^{N}(x_n-\mu)^T\Sigma^{-1}(x_n-\mu)-\frac{N}{2}\ln\det(\Pi)-\frac{1}{2}(\mu-\mu_0)^T\Sigma_0^{-1}(\mu-\mu_0)+const$$

If $\Sigma$ is given, the problem becomes maximum likelihood with regularization (analogous to linear regression with regularization) :

$$\ln[P(X\mid\mu,\Sigma)\Pr(\mu,\Sigma)] \quad = \quad -\frac{1}{2}\sum_{n=1}^{N}(x_n-\mu)^T\Sigma^{-1}(x_n-\mu)-\frac{1}{2}(\mu-\mu_0)^T\Sigma_0^{-1}(\mu-\mu_0)+const* \qquad\qquad \textit{moving known terms into const}$$

Differentiate log likelihood above wrt $\mu$ and setting it to zero, we have :

$$0 \quad = \quad -\frac{1}{2}\sum_{n=1}^{N}(-2)\Sigma^{-1}(x_n-\mu)^T-\frac{1}{2}(2)\Sigma_0^{-1}(\mu-\mu_0)^T$$

$$0 \quad = \quad \sum_{n=1}^{N}(x_n-\mu)\Sigma^{-T}-(\mu-\mu_0)\Sigma_0^{-T} \qquad\qquad \textit{transpose on both sides}$$

$$\mu(N\Sigma^{-T}+\Sigma_0^{-T}) \quad = \quad (\sum_{n=1}^{N}x_n)\Sigma^{-T}+\mu_0\Sigma_0^{-T}$$

$$\mu \quad = \quad (N\Sigma^{-T}+\Sigma_0^{-T})^{-1}(\sum_{n=1}^{N}x_n)\Sigma^{-T}+\mu_0\Sigma_0^{-T}$$

$$= \quad \textit{weighted average between given } \mu_0 \textit{ and mean of data}$$

$$\textit{weight} \quad \propto \quad \textit{inverse of covariance}$$

## 1.5 Linear Regression

*Please read ASM's AX=B documents for details.*

Comparison between *ML* for mean estimation and *ML* for *AX=B* :

|  | Gaussian paramter estimation | Linear regression |
|---|---|---|
| *A* | *no* | $N\times M$ |
| *B* | *no* | $N\times 1$ |
| *X* | $N\times M$ | $M\times 1$ |
| *variable* | $X_n$ *i.e.* $\Re^M$ *row vector* | $(AX-B)_n$ *i.e.* $\Re^1$ *scalar* |
| $\mu$ | *mean vector of* $X_n$ | *mean vector of* $(AX-B)_n$, *i.e. all zero ideally* |
| $\Sigma$ | *covariance between M axes* | *covariance between N data* |

## 2.1 Logistic Regression

Logistic regression is a maximum likelihood classifier, it is an extension from regression to classification, via treating the regression function as a separating plane. Given $N$ *i.i.d.* labelled data points $\{x_n \in \mathfrak{R}^M, y_n \in [0,1], n \in [1,N]\}$, the objective is to find separating plane $f(x;\theta) = 0$ so that $s(f(x;\theta)) \in \mathfrak{R}^1$ returns a scalar value above *0.5* for class *1* and a scalar value below *0.5* for class *0*, where $s(z)$ is known as sigmoid function or logistic function, which tends to *1* as $z \to \infty$ and tends to *0* as $z \to -\infty$. Sigmoid function is selected because of its mathematical convenience in find first order and second order derivatives :

$$
\begin{aligned}
s(z) &= (1+e^{-z})^{-1} \\
s'(z) &= (1+e^{-z})^{-2}e^{-z} \\
&= \frac{1}{1+e^{-z}}\frac{e^{-z}}{1+e^{-z}} \\
&= s(z)(1-s(z)) \\
s''(z) &= \partial_z s(z)(1-s(z)) \\
&= s(z)(1-s(z))^2 - s^2(z)(1-s(z)) \\
&= s(z) - 2s^2(z) + s^3(z)
\end{aligned}
$$

### 2.1a Define likelihood
Let's define likelihood for one data point :

$$
\begin{aligned}
P(x=x_n, y=1\,|\,\theta) &= s(f(x_n\,|\,\theta)) \\
P(x=x_n, y=0\,|\,\theta) &= 1-s(f(x_n\,|\,\theta)) \\
\text{or} \qquad P(x, y\,|\,\theta) &= [s(f(x\,|\,\theta))]^y [1-s(f(x\,|\,\theta))]^{1-y}
\end{aligned}
$$

Let's define likelihood for one dataset :

$$
\begin{aligned}
L(\theta) &= \prod_{n=1}^{N} P(x_n, y_n\,|\,\theta) \\
&= \prod_{n=1}^{N} [s(f(x_n\,|\,\theta))]^{y_n} [1-s(f(x_n\,|\,\theta))]^{1-y_n} \\
\ln L(\theta) &= \sum_{n=1}^{N} y_n \ln[s(f(x_n\,|\,\theta))] + (1-y_n)\ln[1-s(f(x_n\,|\,\theta))] \qquad \textcolor{red}{\textit{looks like Kelly criteria, which has E[], just a coincidence}}
\end{aligned}
$$

### 2.1b Maximise likelihood
Logistic regression is maximizing likelihood by gradient ascent.

$$
\begin{aligned}
\partial_\theta \ln L(\theta) &= \sum_{n=1}^{N} \frac{y_n}{s(f(x_n\,|\,\theta))}\partial_\theta s(f(x_n\,|\,\theta)) - \frac{1-y_n}{1-s(f(x_n\,|\,\theta))}\partial_\theta s(f(x_n\,|\,\theta)) \\
&= \sum_{n=1}^{N}\left[\frac{y_n}{s(f(x_n\,|\,\theta))} - \frac{1-y_n}{1-s(f(x_n\,|\,\theta))}\right]\partial_f s(f) \times \partial_\theta f(x_n\,|\,\theta) \\
&= \sum_{n=1}^{N}\left[\frac{y_n}{s(f(x_n\,|\,\theta))} - \frac{1-y_n}{1-s(f(x_n\,|\,\theta))}\right]s(f(x_n\,|\,\theta))(1-s(f(x_n\,|\,\theta))) \times \partial_\theta f(x_n\,|\,\theta) \\
&= \sum_{n=1}^{N}\left[y_n(1-s(f(x_n\,|\,\theta))) - (1-y_n)s(f(x_n\,|\,\theta))\right] \times \partial_\theta f(x_n\,|\,\theta) \\
&= \sum_{n=1}^{N}\left[y_n - s(f(x_n\,|\,\theta))\right] \times \partial_\theta f(x_n\,|\,\theta) \\
&= J_F^T E
\end{aligned}
$$

$$
\begin{aligned}
\text{where} \qquad J_F[n,m] &= \partial_{\theta_m} f(x_n\,|\,\theta) &\qquad &\textit{$J_F$ is N×M Jacobian matrix} \\
E[n] &= y_n - s(f(x_n\,|\,\theta)) &\qquad &\textit{E is N×1 error col matrix}
\end{aligned}
$$

### 2.1c Linear case
If we pick linear separating plane $f(x;\theta) = x\theta = 0$ (recall that we can always append each vector $x_n$ by *1* and extend to $\mathfrak{R}^{M+1}$, assuming data vectors in $X$ are row vectors, parameter vector $\theta$ is column vector, i.e. *ASM* notation : *AX+B*), we have the first implementation method, gradient ascent with size step $\alpha$ :

$$
\begin{aligned}
\partial_\theta \ln L(\theta) &= \sum_{n=1}^{N}\left[y_n - s(f(x_n\,|\,\theta))\right]x_n^T \\
\Rightarrow \qquad \theta_t &= \theta_{t-1} + \alpha J_F^T E \\
&= \theta_{t-1} + \alpha X^T E
\end{aligned}
$$

## 2.2 Expectation Maximization

Suppose an observable random vector $X \in \mathfrak{R}^{M_1}$ and a non-observable random vector $Z \in \mathfrak{R}^{M_2}$ are drawn from parametric distribution $P(X,Z \mid \theta)$ with unknown parameter vector $\theta$, expectation maximization is a 2-step iterative approach that estimates parameter from one given realization of $X$ using maximum likelihood. Main difference between *ML* and *EM* is that the latter needs to handle <u>latent space</u>, which is the space spanned by hidden data. Among the $M_1+M_2$ dimensions of one realization, $M_2$ are <u>latent variables</u>. Latent variables cannot be directly measured, but they can be inferred from observable variables. For instance, in economics, quality of life and happiness are latent, they can be inferred from wealth, employment & consumption. EM algorithm can be summarised as :

$$\theta_{t+1} \quad = \quad \arg\max_{\theta} E_{Z|X,\theta_t}[\ln P(X,Z \mid \theta)] \qquad \text{where X and Z are observable and non-observable part of one realization set}$$

$$= \quad \arg\max_{\theta} \sum_{Z} [P(Z \mid X,\theta_t)\ln P(X,Z \mid \theta)] \qquad \text{where P is joint distribution of one realization set}$$

$$\theta_{opt} \quad = \quad \lim_{t \to \infty} \theta_t$$

The implementation can be broken down into E-step and M-step. E-step estimates latent variables, M-step estimates parameters.

*E-step :*      *update expectation by finding $P(Z \mid X,\theta_k)$ given estimated parameters from previous iteration t*

*M-step :*      *maximize expectation with respect to $\theta$ using various numerical methods (analytic solution for Gaussian mixture model)*

### 2.2a Define convexity - three approaches

#### Convexity definition 1

First of all, we need to define convex function and derive its properties. A function is said to be convex on interval *[L,U]* if :

$$f(\lambda x_1 + (1-\lambda)x_2) \quad \leq \quad \lambda f(x_1)+(1-\lambda)f(x_2) \qquad \forall x_1, x_2 \in [L,U], \lambda \in [0,1]$$

#### Convexity definition 2

If a function is twice differentiable on interval *[L,U]* and $f''(x) \geq 0$, $\forall x \in [L,U]$, then the function is convex on interval *[L,U]*, let's prove.

$$x = \lambda x_1 + (1-\lambda)x_2 \qquad \text{hence we have } L \leq x_1 \leq x \leq x_2 \leq H$$

$$\lambda x + (1-\lambda)x = \lambda x_1 + (1-\lambda)x_2$$

$$\lambda(x-x_1) = (1-\lambda)(x_2-x)$$

$$\lambda f'(s_1)(x-x_1) \leq (1-\lambda)f'(s_2)(x_2-x) \qquad \exists s_1 \in [x_1,x] \text{ and } \exists s_2 \in [x,x_2], \text{ since } f'' \geq 0, f' \text{ is increasing, so } f'(s_1) \leq f'(s_2)$$

$$\lambda(f(x)-f(x_1)) \leq (1-\lambda)(f(x_2)-f(x)) \qquad \text{apply mean value thm to LHS : there exists tangent at } s_1 \text{ parallel to secant } [x_1,x]$$

$$\text{apply mean value thm to RHS : there exists tangent at } s_2 \text{ parallel to secant } [x,x_2]$$

$$\lambda f(x)+(1-\lambda)f(x) \leq \lambda f(x_1)+(1-\lambda)f(x_2) \qquad \text{i.e. f is convex within } [L,U]$$

#### Convexity definition 3

We now generalize convexity definition to Jensen inequality, given function *f* being convex on *[L,U]*, then :

$$f(\sum_n \lambda_n x_n) \quad \leq \quad \sum_n \lambda_n f(x_n) \qquad \text{where } x_n \in [L,U], \forall n \text{ and } \lambda_n \geq 0, \forall n \text{ such that } \sum_n \lambda_n = 1$$

The proof is done by mathematical induction. The *N=1* case is trival. Here is case *N* :

$$f(\sum_{n=1}^{N} \lambda_n x_n) = f(\lambda_N x_N + \sum_{n=1}^{N-1} \lambda_n x_n)$$

$$= f(\lambda_N x_N + (1-\lambda_N)\sum_{n=1}^{N-1}(1-\lambda_N)^{-1}\lambda_n x_n)$$

$$\leq \lambda_N f(x_N) + (1-\lambda_N)f(\sum_{n=1}^{N-1}(1-\lambda_N)^{-1}\lambda_n x_n) \qquad \text{apply definition of convexity}$$

$$\leq \lambda_N f(x_N) + (1-\lambda_N)\sum_{n=1}^{N-1}(1-\lambda_N)^{-1}\lambda_n f(x_n) \qquad \text{apply case N-1}$$

$$= \lambda_N f(x_N) + \sum_{n=1}^{N-1}\lambda_n f(x_n)$$

$$= \sum_{n=1}^{N}\lambda_n f(x_n)$$

## 2.2b Convexity of log

Show that *ln(x)* is strictly concave on *(0,∞)* using 2nd order differentiation :

$$\partial_x(-\ln x) \quad = \quad -1/x$$

$$\partial_{xx}(-\ln x) \quad = \quad 1/x^2 \quad > \quad 0 \qquad\qquad \forall x \in (0,\infty)$$

With Jensen inequality, we can derive *"log of sum is greater than sum of log"* which is useful for derivation of EM algorithm.

$$\ln(\textstyle\sum_{n=1}^{N}\lambda_n x_n) \quad \geq \quad \sum_{n=1}^{N}\lambda_n \ln(x_n)$$
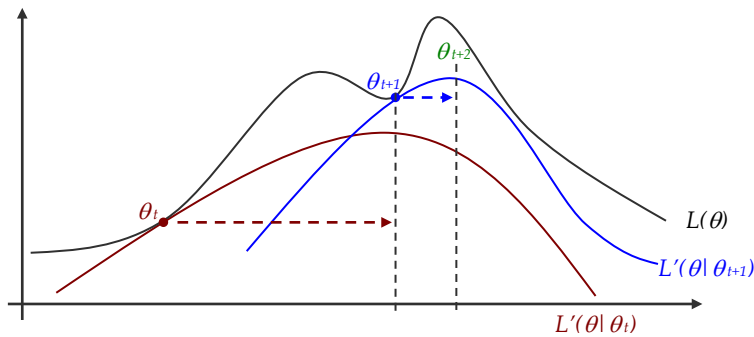
## 2.2c Define likelihood lower bound

Let's derive *EM* algorithm and show its convergence. In general for *EM*, we dont assume *X* to be a set of *N* independent data points. The idea is to derive the likelihood of parameter $\theta$ by incorporating latent space information. Instead of maximising this likelihood, we maximise the lower-bound of the likelihood, as the lower bound is expressed in terms of the likelihood of parameter $\theta_t$ found in iteration *t*, this results in an iterative algorithm that ensures a monotonic increase in likelihood per iteration.

$$
\begin{aligned}
L(\theta) \quad &= \quad \ln P(X \mid \theta) \\
&= \quad \ln \sum_Z P(X,Z \mid \theta) && \textit{step 1 : construct M step probability} \\
&= \quad \ln \sum_Z P(Z \mid X,\theta_t)\frac{P(X,Z \mid \theta)}{P(Z \mid X,\theta_t)} && \textit{step 2 : construct E step probability} \\
&= \quad \ln\left[ \sum_Z P(Z \mid X,\theta_t)\frac{P(X,Z \mid \theta)}{P(Z \mid X,\theta_t)P(X \mid \theta_t)}\right]P(X \mid \theta_t) && \textit{step 3 : construct likelihood of previous iteration } P(X \mid \theta_t) \\
&= \quad \ln\left[ \sum_Z P(Z \mid X,\theta_t)\frac{P(X,Z \mid \theta)}{P(Z \mid X,\theta_t)P(X \mid \theta_t)}\right]+\underbrace{\ln P(X \mid \theta_t)}_{L(\theta_t)} \\
&= \quad \ln\left[ \sum_Z P(Z \mid X,\theta_t)\frac{P(X,Z \mid \theta)}{P(X,Z \mid \theta_t)}\right]+L(\theta_t) \\
&\geq \quad \sum_Z P(Z \mid X,\theta_t)\ln\frac{P(X,Z \mid \theta)}{P(X,Z \mid \theta_t)}+L(\theta_t) && \textit{log of sum greater than sum of log} \\
&\equiv \quad L'(\theta \mid \theta_t) && \textit{define L' as lower bound of log likelihood}
\end{aligned}
$$

We define *L′(θ|θ$_t$)* as the lower bound *L(θ)*. We can show that *L(θ)* and *L′(θ|θ$_t$)* meet at one point :

$$
\begin{aligned}
L'(\theta_t \mid \theta_t) \quad &= \quad \sum_Z P(Z \mid X,\theta_t)\ln\frac{P(X,Z \mid \theta_t)}{P(X,Z \mid \theta_t)}+L(\theta_t) \\
&= \quad \sum_Z P(Z \mid X,\theta_t)\ln 1 + L(\theta_t) \\
&= \quad L(\theta_t)
\end{aligned}
$$

Hence by approximating $L(\theta)$ using its lower bound $L'(\theta|\theta_t)$, which touch at $\theta_t$, we can guarantee an improvement in likelihood by maximizing the lower bound $L'(\theta|\theta_t)$, which leads to EM algorithm.

- $L(\theta) \geq L'(\theta \mid \theta_t)$
- $L(\theta_t) = L'(\theta_t \mid \theta_t)$
- $L(\theta_{t+1}) \geq \underbrace{L'(\theta_{t+1} \mid \theta_t) > L'(\theta_t \mid \theta_t)}_{} = L(\theta_t)$

*update from L(θt) to L(θt+1) can be done by* max L'(θ|θt)

$$
\begin{aligned}
\theta_{t+1} &= \underset{\theta}{\arg\max}\, L'(\theta \mid \theta_t) \\
&= \underset{\theta}{\arg\max} \sum_Z P(Z \mid X, \theta_t) \ln \frac{P(X, Z \mid \theta)}{P(X, Z \mid \theta_t)} + L(\theta_t) \\
&= \underset{\theta}{\arg\max} \sum_Z P(Z \mid X, \theta_t) \ln P(X, Z \mid \theta) \qquad \text{removing all non } \theta \text{ related terms} \\
&= \underset{\theta}{\arg\max}\, E_{Z \mid X, \theta_t}[\ln P(X, Z \mid \theta)]
\end{aligned}
$$

What is gained by trading the maximization of $L(\theta)$ for the maximization of $L'(\theta|\theta_t)$? The answer is that the latter approach offers a framework to take latent space $Z$ into consideration through the two-step iterative algorithm. Prior knowledge about latent space is introduced via $P(Z|X,\theta_t)$.

EM algorithm is thus consisted of two steps :
- estimation of $P(Z \mid X, \theta_t)$ in *expectation step* which will be used as weight for next step
- estimation of $\theta_{t+1} = \underset{\theta}{\arg\max}\, E_{Z \mid X, \theta_t}[\ln P(X, Z \mid \theta)]$ in ***maximization step*** which maximises joint likelihood of $X$ and $Z$

Please read "The Expectation Maximization Algorithm, a short tutorial" by Sean Borman.

## 2.3 Gaussian Mixture Model

Given $\aleph$ modal-GMM, generate a set of $N$ data $[x_1, x_2, x_3, .., x_N]$ by repeatedly :

- picking one *class label* $Z$ from $z \in [1, \aleph]$ with probability $[\pi_1, \pi_2, \pi_3, ..., \pi_\aleph]$
- randomizing *data vector* $X$ from $x \in \Re^M$ with probability $G(x \mid \mu_z, \Sigma_z)$

Given a set of $N$ data $[x_1, x_2, x_3, .., x_N]$, estimate $\aleph$ modal-GMM parameters by *EM* :

- in E step : update the conditional probability of data $x_n$ belonging to a cluster $z$, for all $N$ and all $\aleph$
- in M step : update the estimation of $[\pi_1, \mu_1, \Sigma_1, \pi_2, \mu_2, \Sigma_2, ..., \pi_\aleph, \mu_\aleph, \Sigma_\aleph]$ using conditional probability found in E step as weights

$\aleph$ = *number of clusters*
$N$ = *number of data vectors*
$M$ = *dimension of data vectors*

|  | *rand var set* | *single rand var* | *realized value* | *possible value* |
|---|---|---|---|---|
| *hidden latent space* | $Z$ | $z_n$ | $z$ | $z \in [1,2,3,...,\aleph]$ |
| *observable data space* | $X$ | $x_n$ | $x$ | $x \in \Re^M$ |
| *parameter vector* | $\theta$ | $\begin{bmatrix} \pi_1, \pi_2,...,\pi_\aleph \\ \mu_1, \mu_2,...,\mu_\aleph \\ \Sigma_1, \Sigma_2,...,\Sigma_\aleph \end{bmatrix}$ | | $\begin{bmatrix} \pi_1 + \pi_2 +...+ \pi_\aleph = 1 \\ \mu_1, \mu_2,...,\mu_\aleph \in \Re^M \\ \Sigma_1, \Sigma_2,...,\Sigma_\aleph \in \Re^{M^2} \end{bmatrix}$ |

### 2.3a Define pdf

The complete *pdf* (including latent space) of a sample $X_n$ is :

$$
\begin{aligned}
P(X_n = x, Z_n = z \mid \theta) &= [\pi_1 G(x \mid \mu_1, \Sigma_1)]^{\delta(z=1)} \times [\pi_2 G(x \mid \mu_2, \Sigma_2)]^{\delta(z=2)} \times...\times [\pi_\aleph G(x \mid \mu_\aleph, \Sigma_\aleph)]^{\delta(z=\aleph)} \\
&= \pi_z G(x \mid \mu_z, \Sigma_z)
\end{aligned} \tag{1}
$$

The incomplete *pdf* (excluding latent space) of a sample $X_n$ is :

$$
\begin{aligned}
P(X_n = x \mid \theta) &= \sum_{z=1}^{\aleph} P(X_n = x, Z_n = z \mid \theta) \\
&= \sum_{z=1}^{\aleph} \pi_z G(x \mid \mu_z, \Sigma_z)
\end{aligned} \tag{2}
$$

The *pdf* in latent space is :

$$
\begin{aligned}
P(Z_n = z \mid \theta) &= \int_{-\infty}^{+\infty} P(X_n = x, Z_n = z \mid \theta) dx \\
&= \int_{-\infty}^{+\infty} \pi_z G(x \mid \mu_z, \Sigma_z) dx \\
&= \pi_z \int_{-\infty}^{+\infty} G(x \mid \mu_z, \Sigma_z) dx \\
&= \pi_z
\end{aligned} \tag{3}
$$

### 2.3b E-step

*E-step* is to calculate probability of class label $z_n$ conditional on data vector $x_n$ and current parameter $\theta_t$. To facilitate the process :

$$
\begin{aligned}
\text{We define} \quad w_{nz} &= P(Z_n = z \mid X_n = x, \theta_t = \theta) \qquad \text{which is the prob inside } E_{Z \mid X, \theta_t}[\ln P(X, Z \mid \theta)] \\
&= \frac{P(X_n = x, Z_n = z \mid \theta_t = \theta)}{P(X_n = x \mid \theta_t = \theta)} \\
&= \frac{\pi_z G(x \mid \mu_z, \Sigma_z)}{\sum_{z'=1}^{\aleph} \pi_{z'} G(x \mid \mu_{z'}, \Sigma_{z'})} \qquad \text{substitute (1) and (2)}
\end{aligned}
$$

## 2.3c M-step

*M-step* is to update parameter estimation by maximizing expected log likelihood *under the measure* found in *E-step* :

$$\theta_{t+1} = \underset{\theta}{\arg\max}\, E_{Z|X,\theta_t}[\ln P(X,Z\,|\,\theta)]$$

$$= \underset{\theta}{\arg\max}\, E_{Z|X,\theta_t}[\ln \textstyle\prod_{n=1}^{N} \pi_{z_n} G(x_n\,|\,\mu_{z_n},\Sigma_{z_n})] \qquad\qquad substitute\ (1)$$

$$= \underset{\theta}{\arg\max}\, E_{Z|X,\theta_t}[\textstyle\sum_{n=1}^{N}\ln(\pi_{z_n}G(x_n\,|\,\mu_{z_n},\Sigma_{z_n}))] \qquad\qquad log\ of\ product \to sum\ of\ log$$

$$= \underset{\theta}{\arg\max}\, \textstyle\sum_{n=1}^{N} E_{Z|X,\theta_t}[\ln(\pi_{z_n}G(x_n\,|\,\mu_{z_n},\Sigma_{z_n}))] \qquad\qquad expectation\ of\ sum \to sum\ of\ expectation$$

$$= \underset{\theta}{\arg\max}\, \textstyle\sum_{n=1}^{N}\sum_{z=1}^{\aleph} w_{nz}\ln(\pi_z G(x_n\,|\,\mu_z,\Sigma_z))] \qquad\qquad expectation\ random\ variable\ z_n\ replaced\ by\ realized\ value\ z$$

$$= \underset{\theta}{\arg\max}\, \textstyle\sum_{n=1}^{N}\sum_{z=1}^{\aleph}\left[ w_{nz}\left( \ln \pi_z - \frac{1}{2}rank(\Sigma_z)\ln(2\pi) - \frac{1}{2}\ln \det(\Sigma_z) - \frac{1}{2}(x_n - \mu_z)\Sigma_z^{-1}(x_n - \mu_z)^T \right) \right] \qquad (4)$$

$$= linear\ wrt\ each\ parameter\ in\ \{\pi_z\}$$
$$\quad quad\ wrt\ each\ parameter\ in\ \{\mu_z, \Sigma_z\}$$

By taking derivative *wrt* one of the mean value $\mu_z$ and setting zero, we have :

$$0 = \frac{\partial}{\partial\mu_z}\textstyle\sum_{n=1}^{N}\sum_{z'=1}^{\aleph}\left[ w_{nz'}\left( -\frac{1}{2}(x_n - \mu_{z'})\Sigma_{z'}^{-1}(x_n - \mu_{z'})^T \right) \right]$$

$$0 = \textstyle\sum_{n=1}^{N} w_{nz}(x_n - \mu_z)$$

$$\mu_{z,t+1} = \frac{\sum_{n=1}^{N} w_{nz} x_n}{\sum_{n=1}^{N} w_{nz}}$$

By taking derivative *wrt* one of the mean value $\Sigma_z$ and setting zero, *without proof* we have :
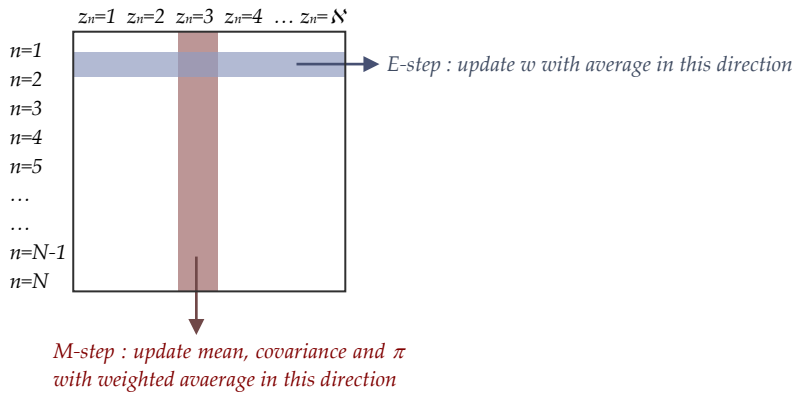
$$\Sigma_{z,t+1} = \frac{\sum_{n=1}^{N} w_{nz}(x_n - \mu_{z,t})^T(x_n - \mu_{z,t})}{\sum_{n=1}^{N} w_{nz}}$$

Equation *(4)* is linear *wrt* $\pi_z$ under constraint $\Sigma\pi_z = 1$, it should be solved using Lagrangian. Without proof, we have :

$$\pi_{z,t+1} = \frac{\sum_{n=1}^{N} w_{nz}}{\sum_{n=1}^{N}\sum_{z'=1}^{\aleph} w_{nz'}}$$

## 2.3d Algorithm using weight matrix

It facilitates the algorithm by defining weight matrix : $w_{n,z} = P(Z_n = z\,|\,X_n = x, \theta_t = \theta)$ .



$z_n=1\ \ z_n=2\ \ z_n=3\ \ z_n=4\ \ \ldots\ z_n=\aleph$

$n=1$
$n=2$ → *E-step : update w with average in this direction*
$n=3$
$n=4$
$n=5$
…
…
$n=N-1$
$n=N$

*M-step : update mean, covariance and $\pi$*
*with weighted avaerage in this direction*

13

## 2.4 K Mean Clustering

Given a set of $N$ observations $\{x_n \in \mathfrak{R}^M, n \in [1,N]\}$ with labels $z_n \in [1, \aleph]$, the objective is to minimise error :

$$\min_{\mu}\left[\sum_n \min_z |x_n - \mu_z|^2\right]$$

The implementation can be broken down into assignment step and update step.

*assign step :*   $z_n = \min_z |x_n - \mu_z|^2$   *E-step in EM : assign each data to closest cluster mean*

*update step :*   $\mu_z = \dfrac{\sum_n \delta(z_n = z)x_n}{\sum_n \delta(z_n = z)}$   *M-step in EM : update cluster mean with assigned data*

K mean clustering can be considered as a special case of Gaussian mixture *EM* having :
- minimising Euclidean distance instead of maximising likelihood
- covariance being removed

## 2.5 K Nearest Neighbour

Given $N$ training vectors having $M$ attributes together with class labels $x_n \in \mathfrak{R}^M$ $y_n \in [1, \aleph]$ $\forall n \in [1,N]$, there are $\aleph$ possible classes. *KNN* classifies test vector $x \in \mathfrak{R}^M$ by considering the $K$ nearest neighbours according to Euclidean distance :

$$n_1 = \arg\min_{n \in [1,N]} |x - x_n|^2$$

$$n_2 = \arg\min_{\substack{n \in [1,N] \\ n \neq n_1}} |x - x_n|^2$$

$$n_3 = \arg\min_{\substack{n \in [1,N] \\ n \neq n_1 n_2}} |x - x_n|^2$$

$$n_4 = \arg\min_{\substack{n \in [1,N] \\ n \neq n_1 n_2 n_3}} |x - x_n|^2$$

$$\ldots$$

$$y = \mathrm{mod}\{n_1, n_2, n_3, \ldots, n_K\}$$

*KNN* is a non-parametric classifer without explicit training step.
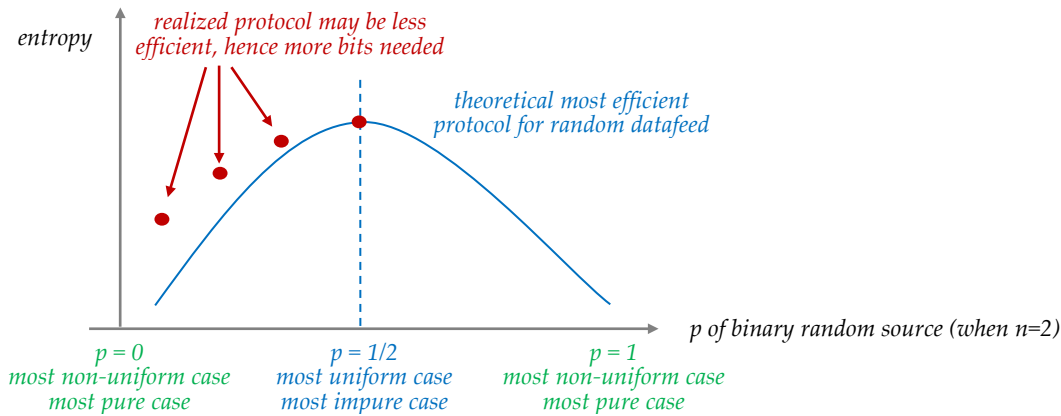
**2.6 KL Divergence** *(based on Entropy)*

*2.6a Define entropy*

Deterministic data source contributes no information (zero entropy), a fair coin flip generates exactly 1 bit information, while unfair coin flip generates less than 1 bit information. Entropy is defined in three ways :

- theoretical minimum number of bits used for serializing stream of datafeed
- measurement of information in number of bits produced by random data source   *(uniform randomness, bigger entropy)*
- measurement of impurity *within **one** output class* produced by separating hyperplane   *(impure classification, bigger entropy)*

$$
\begin{aligned}
entropy \quad &= \quad \sum_n (-p_n \log_2 p_n) \\
&= \quad E_P[-\log_2 p_n]
\end{aligned}
$$



*Examples*

Given a datafeed as a sequence of labels randomly drawn from set *{A,B,C}* with probability *{1/2,1/4,1/4}*, then the optimum protocol should be *{A=0,B=10,C=11}*, thus a stream *{1001100100}* means *{BACAABA}*, entropy and expected length of protocol are :

$$
entropy \quad = \quad \frac{1}{2}\log_2 2^1 + \frac{1}{4}\log_2 2^2 + \frac{1}{4}\log_2 2^2 \qquad = \quad 1.5\ bits
$$

$$
E[length] \quad = \quad \frac{1}{2}length(A) + \frac{1}{4}length(B) + \frac{1}{4}length(C) \qquad = \quad 1.5\ bits
$$

As the expected length of proposed protocol equals to the entropy implying that the protocol is optimum. Lets consider another set *{A,B,C,D}* with probability *{1/2,1/4,1/8,1/8}*, the optimum protocol should be *{A=0,B=10,C=110,D=111}*, let's check :

$$
entropy \quad = \quad \frac{1}{2}\log_2 2^1 + \frac{1}{4}\log_2 2^2 + \frac{1}{8}\log_2 2^3 + \frac{1}{8}\log_2 2^3 \qquad = \quad 1.75\ bits
$$

$$
E[length] \quad = \quad \frac{1}{2}length(A) + \frac{1}{4}length(B) + \frac{1}{8}length(C) + \frac{1}{8}length(D) \quad = \quad 1.75\ bits
$$

Enthopy gives a theoretical minimum, however in practice, the optimum protocol may not be achievable. Consider the set *{A,B,C,D}* with probability *{5/8,1/8,1/8,1/8}*, the best achievable protocol *{A=0,B=10,C=110,D=111}* is no longer optimum :

$$
entropy \quad = \quad \frac{5}{8}\log_2 (2^3/5) + \frac{1}{8}\log_2 2^3 + \frac{1}{8}\log_2 2^3 + \frac{1}{8}\log_2 2^3 \qquad \sim \quad 1.549\ bits
$$

$$
E[length] \quad = \quad \frac{5}{8}length(A) + \frac{1}{8}length(B) + \frac{1}{8}length(C) + \frac{1}{8}length(D) \quad = \quad 1.625\ bits
$$

KL divergence is defined as the number of extra bits required to serialize sample stream from distribution $P$ using optimal protocol for distribution $Q$ relative to using optimal protocol for distribution $P$ itself. Mathematically, it can be denoted as :

*depends on randomness of datafeed (pdf)*
*depends on protocol definition (#bits)*

$$D_{KL}(P \| Q) \quad = \quad \sum_n (-p_n \log_2 q_n) - \sum_n (-p_n \log_2 p_n) \quad = \quad E_P[-\log_2 q_n] - E_P[-\log_2 p_n]$$

$$= \quad \sum_n (-p_n \log_2 \frac{q_n}{p_n}) \quad = \quad E_P\left[-\log_2 \frac{q_n}{p_n}\right]$$

KL divergence in continuous domain is :

$$D_{KL}(P \| Q) \quad = \quad -\int P(X) \ln \frac{Q(X)}{P(X)} dX \quad = \quad E_P\left[-\log_2 \frac{Q(X)}{P(X)}\right]$$

KL divergence is **NOT** the difference in entropy between two distributions :

$$D_{KL}(P \| Q) \quad \neq \quad E_Q[-\log_2 q_n] - E_P[-\log_2 p_n] \qquad \text{which can be positive or negative}$$

KL divergence is always non-negative, it is zero when $P$ and $Q$ are the same distribution. Here is the proof by Jensen inequality.

$$D_{KL}(P \| Q) \quad = \quad \sum_n (-p_n \log_2 \frac{q_n}{p_n})$$

$$\geq \quad -\log_2 \left( \sum_n (p_n \frac{q_n}{p_n}) \right) \qquad \text{since } -ln(x) \text{ is convex, hence } -ln(\Sigma w_n x_n) \leq -\Sigma w_n ln(x_n)$$

$$= \quad -\log_2 \left( \sum_n q_n \right)$$

$$= \quad -\log_2 1$$

$$= \quad 0 \qquad \text{hence when 2 distributions are equivalent, they have zero KL div}$$

KL divergence offers another approach for model parameter estimation. Given a set of $N$ data points $x_n \in \Re^m \ \forall n \in [1,N]$, we contruct an empirical data distribution $P(X|\{x_n\})$, and define parametric distribution model $Q(X|\theta)$, then we can estimate parameter vector $\theta$ by minimizing KL divergence, as it is equivalent to maximization of likelihood $Q(\{x_n\}|\theta)$. Let's prove, the empirical data distribution is simply a summation of delta function locating on data points. In the following, we use continuous version KL divergence.

$$\underset{\theta}{\arg\min}\, D_{KL}(P(X|\{x_n\}) \| Q(X|\theta)) \quad = \quad \underset{\theta}{\arg\min}\, E_P\left[-\log_2 \frac{Q(X)}{P(X)}\right]$$

$$= \quad \underset{\theta}{\arg\max} \int P(X|\{x_n\}) \ln \frac{Q(X|\theta)}{P(X|\{x_n\})} dX$$

$$= \quad \underset{\theta}{\arg\max} \int P(X|\{x_n\}) \ln Q(X|\theta) dX - \int P(X|\{x_n\}) \ln P(X|\{x_n\}) dX$$

$$= \quad \underset{\theta}{\arg\max} \int P(X|\{x_n\}) \ln Q(X|\theta) dX$$

$$= \quad \underset{\theta}{\arg\max} \int \frac{1}{N} \sum_n \delta(x - x_n) \ln Q(X|\theta) dX \qquad \text{since } P(X|\{x_n\}) = \frac{1}{N}\sum_n \delta(x - x_n)$$

$$= \quad \underset{\theta}{\arg\max} \frac{1}{N} \sum_n \ln Q(x_n|\theta)$$

$$= \quad \text{maximum likelihood}$$

KL divergence is closely related to various machine learning algorithms :
- maximum likelihood
- independent component analysis
- generative adversarial network

## 2.7 Decision Tree

### 2.7a Inspection

Decision tree is a tree (not necessarily binary tree) such that each node should either have no child or more than one children. Each node is a separating plane in *ONE* particular dimension of the feature space. On receiving an input feature vector, a node will then forward the vector to one of its child (subtree) according to the separating plane of the node (which is simply a thresholding of that particular dimension in the feature space). The process is recursive and terminated when the input vector reaches one of the leaves, each leaf should be labelled with a class, this is the output of the input feature vector. In short, inspection in decision tree is simply a tree traversal with repeated decision making until a terminal node is reached.

### 2.7b Tree construction

Tree is constructed with :

- each non-terminating node acting as a decision boundary (thresholding in one selected feature)
- each terminating node acting as a class label output
- the tree is parameterized by : hyperplanes at non-terminating nodes and labels at terminating nodes

### 2.7c Tree training

Training is to find all hyperplanes using a greedy algorithm. Starting from the root recursively search the optimal split for the node of interest. For a *K-class* decision tree, let's consider a particular node of interest :

- select one out of *M* features
- select number of children it has, suppose it is *L*, which can be different from the number of class *K*
- select non-overlapping feature range for each of the *L* children
- calculate impurity (entropy) of each child and
- calculate weighted average of entropy, or equivalently, impurity

```cpp
// suppose K is customizable, while L is fixed at 3
double node::average_entropy(const std::vector<data>& training, const std::string& feature_id, double thres_L, double thres_H)
{
    std::vector<data> vecA; // training data fwd to childA
    std::vector<data> vecB; // training data fwd to childB
    std::vector<data> vecC; // training data fwd to childC
    for(const auto& x:training)
    {
        const auto& feature = x.get(feature_id);
        if      (feature < thres_L) vecA.push_back(x);
        else if (feature < thres_H) vecB.push_back(x);
        else                        vecC.push_back(x);
    }
    double HA = entropy(vecA); double wA = vecA.size()/training.size();
    double HB = entropy(vecB); double wB = vecB.size()/training.size();
    double HC = entropy(vecC); double wC = vecC.size()/training.size();
    return wA*HA + wB*HB + wC*HC;
}

double node::entropy(const std::vector<data>& vec)
{
    std::vector<int> counts(K,0); // vector of K zeros
    for(const auto& x:vec) ++counts[x.label];

    double output = 0;
    for(const auto& n:counts)
    {
        double p = n/vec.size();
        output += -p * log2(p);
    }
    return output;
}
```

We have to search `feature_id` together with appropriate `thresholds` to minimize `average_entropy`. Mathematically, we have :

$$\min_{\substack{m\in[1,M] \\ threshold}} \frac{N_A}{N} entropy_A(m,threshold) + \frac{N_B}{N} entropy_B(m,threshold) + \frac{N_C}{N} entropy_C(m,threshold) \qquad where \ N = N_A + N_B + N_C$$

$$entropy_A(m,threshold) = -\sum_k \frac{N_{Ak}}{N_A}\left(\log_2 \frac{N_{Ak}}{N_A}\right) and \ N_{Ak} = num \ of \ data \ labelled \ k \ fwd \ to \ childA \qquad where \ N_A = \sum_k N_{Ak}$$

$$entropy_B(m,threshold) = -\sum_k \frac{N_{Bk}}{N_B}\left(\log_2 \frac{N_{Bk}}{N_B}\right) and \ N_{Bk} = num \ of \ data \ labelled \ k \ fwd \ to \ childB \qquad where \ N_B = \sum_k N_{Bk}$$

$$entropy_C(m,threshold) = -\sum_k \frac{N_{Ck}}{N_C}\left(\log_2 \frac{N_{Ck}}{N_C}\right) and \ N_{Ck} = num \ of \ data \ labelled \ k \ fwd \ to \ childC \qquad where \ N_C = \sum_k N_{Ck}$$

### 3.1 Kalman Filter

*Reference "How a Kalman filter works, in pictures", a Bzarg blog.*

Given Markovian random process with Gaussian noise (where $n$ is time index, $F_n$ and $B_n$ are time-dependent known-dynamics) :

$$X_n \quad = \quad \underbrace{F_n}_{M \times M} \underbrace{X_{n-1}}_{M \times 1} + \underbrace{B_n}_{M \times 1} + \underbrace{\varepsilon(0, Q_n)}_{M \times 1} \qquad\qquad where \quad \underbrace{Q_n}_{M \times M} = E(\varepsilon \varepsilon^T)$$

where $X_{n-1}$ is random. Suppose we keep tracing $X_n$ since beginning of time, we can perform prediction using maximum likelihood :

$$
\begin{aligned}
\hat{X}_n \quad &= \quad \arg\max_X p(X_n = X \mid X_{n-1} = \hat{X}_{n-1}) \\[2mm]
&= \quad \arg\max_X \frac{1}{\sqrt{(2\pi)^M \det(Q_n)}} \exp[-(X - (F_n \hat{X}_{n-1} + B_n))^T Q_n^{-1} (X - (F_n \hat{X}_{n-1} + B_n))/2] \\[2mm]
&= \quad F_n \hat{X}_{n-1} + B_n \qquad\qquad \textit{since ML estimation of Gaussian random variable is just its mean vector}
\end{aligned}
$$

This estimation is theoretically unbiased if the previous estimation is unbiased :

$$
\begin{aligned}
E[X_n - \hat{X}_n] \quad &= \quad E[X_n - (F_n \hat{X}_{n-1} + B_n)] \\[1mm]
&= \quad E[F_n X_{n-1} + B_n + \varepsilon(0, Q_n) - (F_n \hat{X}_{n-1} + B_n)] \\[1mm]
&= \quad E[F_n(X_{n-1} - \hat{X}_{n-1})] + E[\varepsilon(0, Q_n)] \\[1mm]
&= \quad F_n E[X_{n-1} - \hat{X}_{n-1}] \\[1mm]
&= \quad 0 \qquad\qquad \textit{since } E[X_{n-1} - \hat{X}_{n-1}] = 0
\end{aligned}
$$

but the variance of estimation error keeps increasing, as variance is accumulated through time :

$$
\begin{aligned}
V[X_n - \hat{X}_n] \quad &= \quad V[X_n - (F_n \hat{X}_{n-1} + B_n)] \\[1mm]
&= \quad V[F_n X_{n-1} + B_n + \varepsilon(0, Q_n) - (F_n \hat{X}_{n-1} + B_n)] \\[1mm]
&= \quad V[F_n(X_{n-1} - \hat{X}_{n-1})] + V[\varepsilon(0, Q_n)] \\[1mm]
&= \quad V[F_n(X_{n-1} - \hat{X}_{n-1})] + Q_n \\[1mm]
&= \quad E[F_n(X_{n-1} - \hat{X}_{n-1})(F_n(X_{n-1} - \hat{X}_{n-1}))^T] + Q_n \\[1mm]
&= \quad F_n E[(X_{n-1} - \hat{X}_{n-1})(X_{n-1} - \hat{X}_{n-1})^T] F_n^T + Q_n
\end{aligned}
$$

For convenience, we declare $P_n$ as the variance of error, hence we have :

$$P_n \quad = \quad F_n P_{n-1} F_n^T + Q_n \qquad\qquad \textit{accumulation of variance through time}$$

If we are given an observation in measurement domain right after each prediction step, we can then update our current estimation for the sake of the next prediction, again this is done using maximum likelihood, the conversion from state domain to measurement domain is done by a time-dependent transformation $H_n$, where $Z_n$ is the measurement reading, again there is measurement noise :

$$Z_n \quad = \quad \underbrace{\varepsilon'(Z_n^{read}, R_n)}_{M' \times 1} \qquad\qquad where \quad \underbrace{R_n}_{M' \times M'} = E(\varepsilon' \varepsilon'^T)$$

Our prediction in measurement domain is $H_n X_n$, mean and variance error of our estimation in measurement space is :

$$
\begin{aligned}
E[H_n(X_n - \hat{X}_n)] \quad &= \quad H_n E[X_n - \hat{X}_n] \quad = \quad 0 \\[2mm]
V[H_n(X_n - \hat{X}_n)] \quad &= \quad E[H_n(X_n - \hat{X}_n)(H_n(X_n - \hat{X}_n))^T] - E[H_n(X_n - \hat{X}_n)]E[H_n(X_n - \hat{X}_n)]^T \\[1mm]
&= \quad E[H_n(X_n - \hat{X}_n)(X_n - \hat{X}_n)^T H_n^T] - H_n E[X_n - \hat{X}_n](H_n E[X_n - \hat{X}_n])^T \\[1mm]
&= \quad H_n E[(X_n - \hat{X}_n)(X_n - \hat{X}_n)^T] H_n^T \\[1mm]
&= \quad H_n P_n H_n^T
\end{aligned}
$$

In this update step, we maximize the joint probability between :

(1) $\quad p(Z_n = Z \mid X_{n-1} = \hat{X}_{n-1}) \qquad = \quad \varepsilon(H_n \hat{X}_n, H_n P_n H_n^T) \qquad$ *distribution of prediction in measurement domain*

(2) $\quad p(Z_n = Z) \qquad\qquad\qquad\qquad = \quad \varepsilon(Z_n^{read}, R_n) \qquad\qquad$ *distribution of observation in measurement domain*

The joint likelihood of prediction and observation in measurement domain is a product of two Gaussians, which is also a Gaussian.

$$\hat{Z}_n \quad = \quad \arg\max_{Z} \varepsilon_Z(H_n\hat{X}_n, H_nP_nH_n^T)\varepsilon_Z(Z_n^{read}, R_n)$$

$$= \quad \arg\max_{Z} \varepsilon_Z \begin{pmatrix} mean = W_A H_n\hat{X}_n + W_B Z_n^{read} \\ var = ((H_nP_nH_n^T)^{-1} + R_n^{-1})^{-1} \end{pmatrix}$$

*where*
$$\begin{bmatrix} W_A \\ W_B \end{bmatrix} = \begin{bmatrix} ((H_nP_nH_n^T)^{-1} + R_n^{-1})^{-1}(H_nP_nH_n^T)^{-1} \\ ((H_nP_nH_n^T)^{-1} + R_n^{-1})^{-1}R_n^{-1} \end{bmatrix}$$     *These matrices look like those in AX=B doc.*

$$= \begin{bmatrix} R_n(H_nP_nH_n^T + R_n)^{-1} \\ H_nP_nH_n^T(H_nP_nH_n^T + R_n)^{-1} \end{bmatrix}$$     *refer to product of Gaussians*

$$= \begin{bmatrix} I - K \\ K \end{bmatrix}$$     *refer to product of Gaussians, define gain vector* $K = W_B$

*Since*
$$K \quad = \quad W_B$$
$$= \quad I - W_A$$
$$= \quad I - ((H_nP_nH_n^T)^{-1} + R_n^{-1})^{-1}(H_nP_nH_n^T)^{-1}$$
$$(I - K)(H_nP_nH_n^T) \quad = \quad ((H_nP_nH_n^T)^{-1} + R_n^{-1})^{-1}$$     *equation #*

For machine learning, we prefer to express as updating equation with Kalman factor as step size (instead of weighted average) :

$$\hat{Z}_n \quad = \quad W_A H_n\hat{X}_n + W_B Z_n^{read}$$     *MLE as weighted sum*
$$= \quad (I - K)H_n\hat{X}_n + KZ_n^{read}$$
$$= \quad H_n\hat{X}_n + K(Z_n^{read} - H_n\hat{X}_n)$$     *MLE as an updated equation*
$$\tilde{X}_n \quad = \quad H_n^{-1}\hat{Z}_n$$
$$= \quad \hat{X}_n + K'(Z_n^{read} - H_n\hat{X}_n)$$     *where* $K' = P_nH_n^T(H_nP_nH_n^T + R_n)^{-1}$

If we use updated estimate instead of the original estimate, accumulation of variance error in measurement domain becomes :

$$H_n\tilde{P}_nH_n^T \quad = \quad ((H_nP_nH_n^T)^{-1} + R_n^{-1})^{-1}$$
$$= \quad (I - K)(H_nP_nH_n^T)$$     *using equation #*
$$= \quad H_nP_nH_n^T - K(0 - H_nP_nH_n^T)$$
$$\tilde{P}_nH_n^T \quad = \quad P_nH_n^T - K'(0 - H_nP_nH_n^T)$$
$$\tilde{P}_n \quad = \quad P_n - K'(0 - H_nP_n)$$     *variance estimation as an updated equation*

Please note that $\hat{X}_n \neq \tilde{X}_n$, though both are ML estimations, they are best estimates given different informations. With updating step, the prediction step is based on updated estimate and becomes :

$$\hat{X}_n \quad = \quad \arg\max_{X} p(X_n = X \mid X_{n-1} = \hat{X}_{n-1})$$     *prediction without updating*
$$\Rightarrow \quad \hat{X}_n \quad = \quad \arg\max_{X} p(X_n = X \mid X_{n-1} = \tilde{X}_{n-1})$$     *prediction with updating*

Difference in roles of $Q$ in prediction and $R$ in updating : in prediction, we add $X$ with $\varepsilon$, while in updating, we join $HX$ with $\varepsilon$ $(Z,R)$.

$$X_n \quad = \quad F_nX_{n-1} + B_n + \varepsilon(0, Q_n)$$     *which results in sum of variances*
$$H_nX_{n-1} \quad vs \quad \varepsilon(Z_n^{read}, R_n)$$     *which results in inverse variance weighs*

## Kalman filter appendix – comparison with Recursive Least Square (RLS)

We can simulate RLS with a Kalman filter by removing predict step of Kalman filter and setting the dimension of measurement space to *1*. Here is the correspondence between the variables in RLS and Kalman filter :

| | Kalman filter | | Recursive least square | | |
|---|---|---|---|---|---|
| *predict step :* | $\hat{X}_n$ | $= F_n \tilde{X}_{n-1} + B_n$ | *no* | | |
| | $P_n$ | $= F_n P_{n-1} F_n^T + Q_n$ | *no* | | |
| | | | | | |
| *update step :* | $H_n$ | $M' \times M$ | $A^{(n)}{}_n$ | $1 \times M$ | *measurement trasnformation* |
| | $Z_n$ | $M' \times 1$ | $b^{(n)}{}_n$ | $1 \times 1$ | *measurement space (with M'=1)* |
| | $R_n$ | $M' \times M'$ | $\sigma_n{}^2$ | $1 \times 1$ | *measurement error* |
| | | | | | |
| | $K'$ | $= P_n H_n^T (H_n P_n H_n^T + R_n)^{-1}$ | $K'$ | $= \dfrac{P_{n-1} A_n^{(n)T}}{A_n^{(n)} P_{n-1} A_n^{(n)T} + \sigma_n^2}$ | |
| *state estimate* | $\tilde{X}_n$ | $= \hat{X}_n + K'(Z_n^{read} - H_n \hat{X}_{n|})$ | $\tilde{X}_n$ | $= X_{n-1} + K'(b_n^{(n)} - A_n^{(n)} X_{n-1})$ | |
| *state covariance* | $\tilde{P}_n$ | $= P_n + K'(0 - H_n P_n)$ | $\tilde{P}_n$ | $= P_{n-1} + K'(0 - A_n^{(n)} P_{n-1})$ | |

## Kalman filter appendix – multiplication of multivariate Gaussians

When two Gaussians sum together, it becomes a Gaussian as convolution of Gaussians gives a Gaussian. When two Gaussians join together (i.e. joint probability in maximum likelihood), it also becomes a Gaussian as product of Gaussians gives a Gaussian. These two results are consistent as we know that Fourier transform of a Gaussian is a Gaussian, while convolution in time is equivalent to product in frequency domain and vice versa. Please make a summary here :

| | | | |
|---|---|---|---|
| $p(X_1 = x)$ | $=$ | $\varepsilon(\mu_1, \sigma_1)$ | |
| $p(X_2 = x)$ | $=$ | $\varepsilon(\mu_2, \sigma_2)$ | |
| $p(X_1 + X_2 = x)$ | $=$ | $\varepsilon(\mu_1, \sigma_1) \otimes \varepsilon(\mu_2, \sigma_2)$ | *sum of Gaussians* |
| | $=$ | $\varepsilon(\mu_1 + \mu_2, \sqrt{\sigma_1^2 + \sigma_2^2})$ | |
| $p(X_1 = x) p(X_2 = x)$ | $=$ | $\varepsilon(\mu_1, \sigma_1) \times \varepsilon(\mu_2, \sigma_2)$ | *joint of Gaussians* |
| | $=$ | $\varepsilon(w_1 \mu_1 + w_2 \mu_2, (\sigma_1^{-2} + \sigma_2^{-2})^{-1})$ | |
| | $=$ | $\varepsilon(w_1 \mu_1 + w_2 \mu_2, w_1 \sigma_1^2 = w_2 \sigma_2^2)$ | |
| *where* $\quad w_1$ | $=$ | $\sigma_2^2 / (\sigma_1^2 + \sigma_2^2)$ | |
| $\quad w_2$ | $=$ | $\sigma_1^2 / (\sigma_1^2 + \sigma_2^2)$ | |

## *Proof of product of Gaussians*

| | | | |
|---|---|---|---|
| $p_k(x)$ | $=$ | $\dfrac{1}{\sqrt{(2\pi)^M \det(P_k)}} \exp(-\dfrac{1}{2}(X - E_k)^T P_k^{-1}(X - E_k))$ | *where* $k = 1,2$ |
| $q(x)$ | $=$ | $p_1(x) p_2(x) / s_1$ | *where* $s_1 = \int_{-\infty}^{\infty} p_1(x) p_2(x) dx$ *(by note 1)* |
| | $=$ | $\exp(-\dfrac{1}{2}(X - E_1)^T P_1^{-1}(X - E_1) - \dfrac{1}{2}(X - E_2)^T P_2^{-1}(X - E_2)) \Big/ s_2$ | *where* $s_2 = \dfrac{s_1}{(2\pi)^M \sqrt{\det(P_1)\det(P_2)}}$ |
| | $=$ | $\exp(-\dfrac{1}{2}(X^T (P_1^{-1} + P_2^{-1})X - X^T(P_1^{-1}E_1 + P_2^{-1}E_2) - (E_1^T P_1^{-1} + E_2^T P_2^{-1})X)) \Big/ s_3$ | *where* $s_3 = \dfrac{s_2}{\exp(-\dfrac{1}{2}(E_1^T P_1^{-1}E_1 + E_2^T P_2^{-1}E_2))}$ |
| | $=$ | $\exp(-\dfrac{1}{2}(X^T \underbrace{(P_1^{-1} + P_2^{-1})}_{A}X - X^T\underbrace{(P_1^{-1}E_1 + P_2^{-1}E_2)}_{B} - \underbrace{(P_1^{-1}E_1 + P_2^{-1}E_2)}_{B}{}^T X)) \Big/ s_3$ | $P_1$ *and* $P_2$ *are symmetric* |
| | $=$ | $\exp(-\dfrac{1}{2}(X - A^{-1}B)^T A(X - A^{-1}B)) \Big/ s_4$ | *where* $s_4 = s_3 \exp(-\dfrac{1}{2}(B^T A^{-1}B))$ *(by note 2)* |
| | $=$ | $Gaussian(E = A^{-1}B, P = A^{-1})$ | |

The new mean vector and new covariance matrix are :

| | | | |
|---|---|---|---|
| $E$ | $=$ | $(P_1^{-1} + P_2^{-1})^{-1}(P_1^{-1}E_1 + P_2^{-1}E_2)$ | *define* $K = (P_1^{-1} + P_2^{-1})^{-1} P_2^{-1}$ |
| | $=$ | $(I - K)E_1 + KE_2$ | *weighted average of two original means* |
| $P$ | $=$ | $(P_1^{-1} + P_2^{-1})^{-1}$ | *which is like two resistors in parallel* |

The new mean is weighted average of the two constituents' mean. The covariance matrix evolves like adding one new resistor $P_2$ in parallel with the original circuit $P_1$. Please note that $K = W_B$ can be represented in 4 different ways :

$$
\begin{aligned}
K \quad &= \quad (P_1^{-1} + P_2^{-1})^{-1} P_2^{-1} && (1)\\
&= \quad I - (P_1^{-1} + P_2^{-1})^{-1} P_1^{-1} && (2)\\
&= \quad P_1(P_1 + P_2)^{-1} && (3)\\
&= \quad I - P_2(P_1 + P_2)^{-1} && (4)
\end{aligned}
$$

### Note 1
*We omit all the constant terms and combine them into s because :*

given two pdf with unity integral $\qquad\qquad \int_{-\infty}^{\infty} p_1(x)dx = 1 \ \ and \ \ \int_{-\infty}^{\infty} p_2(x)dx = 1$

integral of their product is non unity $\qquad\quad \int_{-\infty}^{\infty} p_1(x)p_2(x)dx \neq 1$

hence product of pdf should be defined as $\qquad q(x) = \dfrac{p_1(x)p_2(x)}{\int_{-\infty}^{\infty} p_1(x)p_2(x)dx}$ *so that* $\int_{-\infty}^{\infty} q(x)dx = 1$

### Note 2
*Given symmetric square matrix A, completing square in matrix form is :*

$$
\begin{aligned}
& X^T A X - X^T B - B^T X \\
=\ & X^T A X - X^T B - B^T X + B^T A^{-1} B - B^T A^{-1} B \\
=\ & X^T \sqrt{A}\sqrt{A} X - X^T B - B^T X + B^T \sqrt{A}^{-1}\sqrt{A}^{-1} B - B^T A^{-1} B && \text{\textit{as we define }} A = \sqrt{A}\sqrt{A} \text{\textit{ , then }} A^{-1} = (\sqrt{A}\sqrt{A})^{-1} = \sqrt{A}^{-1}\sqrt{A}^{-1} \\
=\ & X^T \sqrt{A}^T \sqrt{A} X - X^T B - B^T X + B^T \sqrt{A}^{-1T}\sqrt{A}^{-1} B - B^T A^{-1} B && \text{\textit{besides, we pick a symmetic }} \sqrt{A}^T = \sqrt{A} \\
=\ & (\sqrt{A}X)^T \sqrt{A} X - X^T B - B^T X + (\sqrt{A}^{-1}B)^T \sqrt{A}^{-1} B - B^T A^{-1} B \\
=\ & (\sqrt{A}X - \sqrt{A}^{-1}B)^T (\sqrt{A}X - \sqrt{A}^{-1}B) - B^T A^{-1} B && \text{\textit{using property 2a and 2b below}} \\
=\ & (\sqrt{A}(X - A^{-1}B))^T \sqrt{A}(X - A^{-1}B) - B^T A^{-1} B \\
=\ & (X - A^{-1}B)^T \sqrt{A}^T \sqrt{A}(X - A^{-1}B) - B^T A^{-1} B && \text{\textit{property 2a}} : (\sqrt{A}X)^T \sqrt{A}^{-1} B = X^T \sqrt{A}^T \sqrt{A}^{-1} B = X^T B \\
=\ & (X - A^{-1}B)^T A(X - A^{-1}B) - B^T A^{-1} B && \text{\textit{property 2b}} : (\sqrt{A}^{-1}B)^T \sqrt{A}X = B^T \sqrt{A}^{-1T} \sqrt{A}X = B^T X \\
& && \text{\textit{recall : transpose of inverse equals to inverse of tranpose}}
\end{aligned}
$$

### Note 3
*More about matrix inverse :*

$$
\begin{aligned}
(A^{-1} + B^{-1})^{-1} &= \quad AA^{-1}(A^{-1} + B^{-1})^{-1} B^{-1} B \\
&= \quad A(B(A^{-1} + B^{-1})A)^{-1} B \\
&= \quad A(BA^{-1}A + BB^{-1}A)^{-1} B \\
&= \quad A(A + B)^{-1} B \\
(A^{-1} + B^{-1})^{-1} &= \quad B(A + B)^{-1} A && \text{\textit{follow similar logic}}
\end{aligned}
$$

*Besides, weights by <u>inverse covariance</u> can be written as weights by <u>covariance</u> :*

$$
\begin{aligned}
W_A \quad &= \quad (\Sigma_A^{-1} + \Sigma_B^{-1})^{-1} \Sigma_A^{-1} \\
&= \quad \Sigma_B(\Sigma_A + \Sigma_B)^{-1} \Sigma_A \Sigma_A^{-1} \\
&= \quad \Sigma_B(\Sigma_A + \Sigma_B)^{-1} && \text{\textit{property 3d}} \\
W_B \quad &= \quad (\Sigma_A^{-1} + \Sigma_B^{-1})^{-1} \Sigma_B^{-1} \\
&= \quad \Sigma_A(\Sigma_A + \Sigma_B)^{-1} \Sigma_B \Sigma_B^{-1} \\
&= \quad \Sigma_A(\Sigma_A + \Sigma_B)^{-1} && \text{\textit{property 3c}} \\
W_A \quad &= \quad (\Sigma_A^{-1} + \Sigma_B^{-1})^{-1} \Sigma_A^{-1} \\
&= \quad (\Sigma_A^{-1} + \Sigma_B^{-1})^{-1} (\Sigma_A^{-1} + \Sigma_B^{-1} - \Sigma_B^{-1}) \\
&= \quad (\Sigma_A^{-1} + \Sigma_B^{-1})^{-1} (\Sigma_A^{-1} + \Sigma_B^{-1}) - (\Sigma_A^{-1} + \Sigma_B^{-1})^{-1} \Sigma_B^{-1} \\
&= \quad I - (\Sigma_A^{-1} + \Sigma_B^{-1})^{-1} \Sigma_B^{-1} \\
&= \quad I - W_B
\end{aligned}
$$

## 3.2 Principal Component Analysis

Given a $N{\times}M$ row-wise data matrix $X$, its mean row-vector $\mu$ and covariance matix $\Sigma$ respectively are :

| | | | |
|---|---|---|---|
| $\mu$ | $=$ | $L^T X / N$ | *where L is a all-one column-vector* |
| $\Sigma$ | $=$ | $(X-U)^T(X-U)/N$ | *where $U = [\mu;\ \mu;\ \mu;\ \dots\mu]$ is a stack of N $\mu$ forming a N$\times$M matrix* |

Principal component analysis finds the projection column-vector $w$ so that variance of transformed data is maximized, transformed mean vector and covariance matrix (both are *1$\times$1*) can be written in terms of the original mean vector and covariance matrix :

| | | | |
|---|---|---|---|
| $\mu'$ | $=$ | $L^T(Xw)/N$ | *we denote <u>data x as row-vector</u> and <u>parameter w as column-vector</u>* |
| | $=$ | $\mu w$ | |
| $\Sigma'$ | $=$ | $(Xw-U')^T(Xw-U')/N$ | *where $U' = [\mu';\ \mu';\ \mu';\ \dots\mu']$ is a stack of N $\mu'$ forming a N$\times$M matrix* |
| | $=$ | $(Xw-Uw)^T(Xw-Uw)/N$ | |
| | $=$ | $w^T(X-U)^T(X-U)w/N$ | |
| | $=$ | $w^T\Sigma w$ | *note that $\mu$ is 1$\times$M, $\Sigma$ is M$\times$M, while both $\mu'$ and $\Sigma'$ are 1$\times$1* |

In order to maximize $\Sigma'$ respect to column-vector $w$ under the constraint $w^Tw = 1$, we set up a Lagrangian, take the gradient and set it zero (gradient is also a column-vector like $w$), we have :

| | | | | |
|---|---|---|---|---|
| $L$ | $=$ | $\Sigma'-\lambda(w^Tw-1)$ | | |
| | $=$ | $w^T\Sigma w - \lambda(w^Tw-1)$ | | |
| $\nabla_w L$ | $=$ | $2\Sigma w - 2\lambda w$ | $=$ | $0$ |
| $\Sigma w$ | $=$ | $\lambda w$ | | |

With eigen decomposition, eigen vector gives the desired projection vector, while eigen value gives the Lagrange multiplier. Yet in practice, we wish to obtain all the principal components $W$, which has column-vectors representing orthogonal basis. As covariance $\Sigma'$ is no longer a scalar, we need some way to define a scalar objective : the answer is trace. Trace of transformed covariance is sum of diagonal variances, while ignoring off-diagonal covariances :

| | | | |
|---|---|---|---|
| $tr(\Sigma')$ | $=$ | $\sum_{n=1}^{N}(\Sigma')_{nn}$ | |
| *and* $\Sigma'$ | $=$ | $W^T\Sigma W$ | *where W is column-wise projection basis, both W and $\Sigma'$ are M$\times$M* |

We are going to prove two properties (1) eigen basis can diagonalize covariance matrix, i.e. off-diagonal covariances of transformed data are zeros, (2) trace of covariance $tr(\Sigma') = tr(W^T\Sigma W)$ is maximized when <u>columns of W</u> are set equal to eigen vectors of original covariance $\Sigma$. First of all, the proof for diagonalization is straight forward, suppose the eigen decomposition of $\Sigma$ is :

| | | | |
|---|---|---|---|
| $\Sigma$ | $=$ | $Q^T\Lambda Q$ | *where Q is row-wise eigen basis, hence $QQ^T{=}I$ and $\Lambda_{nm} = 0\ \forall n{=}m$* |
| $\Rightarrow\quad \Sigma'$ | $=$ | $W^T\Sigma W$ | |
| | $=$ | $W^TQ^T\Lambda QW$ | |
| | $=$ | $(QW)^T\Lambda(QW)$ | |
| $\Sigma'\vert_{W=Q^T}$ | $=$ | $\Lambda$ | *when columns of W are eigen vectors of $\Sigma$, i.e. $W = Q^T$ and $QW = QQ^T = I$* |

Secondly, we have to show that trace of $\Sigma'\vert_{W=QT}$ is greater than trace of $\Sigma'$ for all other $W$. It involves the trace trick :

| | | | |
|---|---|---|---|
| $tr(\Sigma')$ | $=$ | $tr((QW)^T\Lambda(QW))$ | |
| | $=$ | $tr(\Lambda(QW)(QW)^T)$ | *trace is invariant to cyclic permutation* |
| | $\leq$ | $tr(\Lambda)tr((QW)(QW)^T)$ | *prove it later using Cauchy Schwarz inequality (see remark after LDA)* |
| | $=$ | $\sum_{n=1}^{N}\Lambda_{nn}tr(QWW^TQ^T)$ | |
| | $=$ | $\sum_{n=1}^{N}\Lambda_{nn}tr(WW^TQ^TQ)$ | *trace is invariant to cyclic permutation* |
| | $=$ | $\sum_{n=1}^{N}\Lambda_{nn}tr(WW^T)$ | *since $QQ^T = Q^TQ = I$ as it is eigen basis* |
| | $=$ | $\sum_{n=1}^{N}\Lambda_{nn}tr(I)$ | *since $WW^T = W^TW = I$ as it is projection basis (constraint for optimization)* |
| | $=$ | $\sum_{n=1}^{N}\Lambda_{nn}\times rank$ | |

Therefore, our objective function is bounded above by the sum of eigen values scaled by rank of data, equality holds when $QW = I$, that is when $W = Q^T$ (columns of W equal to rows of Q). In short eigen basis is the projection basis that maximise covariance trace.

## 3.3 Linear Discriminant Analysis

If logistic regression is regarded as a supervised linear classifier, then linear discriminant analysis can be regarded as unsupervised linear classifier, the former is derived from maximum likelihood, while the latter is derived from between-class against within-class variance ratio, which can be considered as an extension from *PCA*. Lets consider two-classes problem, given a $N_1{\times}M$ row-wise data matrix $X_1$ for class *1* and a $N_2{\times}M$ row-wise data matrix $X_2$ for class 2, the mean row-vectors $\mu_1\,\mu_2$ and covariance matrices $\Sigma_1\,\Sigma_2$ are :

$$\mu_k \quad = \quad L_k^T X_k / N_k \qquad\qquad\qquad\qquad \textit{where } L_k \textit{ is a all-one column-vector}$$

$$\Sigma_k \quad = \quad (X_k - U_k)^T (X_k - U_k)/N_k \qquad \textit{where } U_k = [\mu_k; \mu_k; \mu_k; \dots \mu_k] \textit{ is a stack of } N_k\, \mu_k \textit{ forming a } N_k{\times}M \textit{ matrix}$$

$$k \quad = \quad 1,2$$

Between class and within class covariance matrix are defined as the following (both matrices are $M{\times}M$) :

$$\Sigma_{between} \quad = \quad \sum_{k=1}^{2}(\mu_k - \mu)^T(\mu_k - \mu) \qquad \textit{where } \mu = LX/N \textit{ and } X = [X_1; X_2],\ N = N_1 + N_2$$

$$\Sigma_{within} \quad = \quad \sum_{k=1}^{2}(X_k - U_k)^T(X_k - U_k)/N_k$$

$$\quad = \quad \sum_{k=1}^{2}\Sigma_k$$

Linear discriminant analysis finds the projection column-vector $w$ such that the between-class and the within-class variance ratio of transformed data is maximized, transformed between-class covariance and transformed within-class covariance are both scalar and they can be written in terms of the original between-class covariance and within-class covariance :

$$\Sigma'_{between} \quad = \quad \sum_{k=1}^{2}(U_k w - U w)^T(U_k w - U w)$$

$$\quad = \quad \sum_{k=1}^{2} w^T(U_k - U)^T(U_k - U)w$$

$$\quad = \quad w^T \Sigma_{between} w$$

$$\Sigma'_{within} \quad = \quad \sum_{k=1}^{2}(X_k w - U_k w)^T(X_k w - U_k w)/N_k$$

$$\quad = \quad \sum_{k=1}^{2} w^T(X_k - U_k)^T(X_k - U_k)w/N_k$$

$$\quad = \quad w^T \Sigma_{within} w$$

Both covariance matrices above are scalar as $w$ is column-vector. Objective of *LDA* is to maximize the following ratio with $w^Tw = 1$ :

$$\max_{w}\frac{\Sigma'_{between}}{\Sigma'_{within}} \quad = \quad \max_{w}\frac{w^T \Sigma_{between} w}{w^T \Sigma_{within} w} \qquad\qquad st\ w^T w = 1$$

$$\quad = \quad \max_{w}\frac{(\alpha w)^T \Sigma_{between}(\alpha w)}{(\alpha w)^T \Sigma_{within}(\alpha w)} \qquad st\ w^T w = 1 \qquad \textit{(as objective is invariant to scaling in } w\textit{)}$$

$$\quad = \quad \max_{w}\frac{(\alpha w)^T \Sigma_{between}(\alpha w)}{w^T(\alpha^T \Sigma_{within}\alpha)w} \qquad st\ w^T w = 1$$

$$\quad = \quad \max_{w}(\alpha w)^T \Sigma_{between}\underbrace{(\alpha w)}_{w'} \qquad st\ w^T w = 1 \qquad \textit{(let's pick } \alpha \textit{ such that } \alpha^T\Sigma_{within}\alpha = I \textit{ , then apply } w^Tw = 1\textit{)}$$

$$\quad = \quad \max_{w'} w'^T \Sigma_{between} w' \qquad st\ (\alpha^{-1}w')^T(\alpha^{-1}w) = 1 \qquad \textit{(put } w' = \alpha w\textit{, hence } w = \alpha^{-1}w'\textit{)}$$

$$\quad = \quad \max_{w'} w'^T \Sigma_{between} w' \qquad st\ w'^T \Sigma_{within} w = 1 \qquad \textit{(since } (\alpha^{-1}w')^T(\alpha^{-1}w) = w'^T(\alpha^{-T}\alpha^{-1})w = w'^T\Sigma_{within}w \textit{ )}$$

Changing notation $w'$ back to $w$, we then set up a Lagrangian, take the gradient and set it zero, we have :

$$L \quad = \quad w^T \Sigma_{between} w - \lambda(w^T \Sigma_{within} w - 1)$$

$$\nabla_w L \quad = \quad 2\Sigma_{between} w - 2\lambda\Sigma_{within} w \quad = \quad 0$$

$$\Sigma_{between} w \quad = \quad \lambda\Sigma_{within} w$$

$$(\Sigma_{within})^{-1}\Sigma_{between} w \quad = \quad \lambda w$$

As matrix $(\Sigma_{within})^{-1}\Sigma_{between}$ is not necessarily symmetric (please see remark for the reason), the above equation isn't an ordinary eigen value decomposition, instead it is known as a generalized eigen problem. In order to deal with ordinary eigen value decomposition, we need to make further simplication by taking square root of $\Sigma_{between}$, that is $\Sigma_{between} = (\Sigma_{between})^{1/2}(\Sigma_{between})^{1/2}$.

$$
\begin{aligned}
(\Sigma_{within})^{-1}\Sigma_{between}^{1/2}\Sigma_{between}^{1/2}w &= \lambda w \\
\Sigma_{between}^{1/2}(\Sigma_{within})^{-1}\Sigma_{between}^{1/2}\Sigma_{between}^{1/2}w &= \lambda \Sigma_{between}^{1/2}w \\
\underbrace{\Sigma_{between}^{1/2}(\Sigma_{within})^{-1}\Sigma_{between}^{1/2}}_{A}w' &= \lambda w' \qquad\qquad by\ putting\ w' = (\Sigma_{between})^{1/2}w
\end{aligned}
$$

Therefore, *LDA* projection can be found by eigen decomposition of matrix $A$ and apply inverse transform $w = (\Sigma_{between})^{-1/2}w'$ on eigen vector $w'$, besides square root of matrix can be found by *SVD*, i.e. if $\Sigma_{between} = U\Lambda U^T$, then we have $(\Sigma_{between})^{1/2} = U\Lambda^{1/2}U^T$. Finally above *LDA* can easily be extended to multiclass by having $K > 2$. For case $K = 2$, most researchers prefer to define between-class covariance as the following :

$$
\Sigma_{betweenII} = (\mu_1 - \mu_2)^T(\mu_1 - \mu_2)
$$

which is equivalent to the original definition, lets prove :

$$
\begin{aligned}
\Sigma_{between} &= \sum_{k=1}^{2}(\mu_k - \mu)^T(\mu_k - \mu) \\
&= (\mu_1 - (r_1\mu_1 + r_2\mu_2))^T(\mu_1 - (r_1\mu_1 + r_2\mu_2)) + (\mu_2 - (r_1\mu_1 + r_2\mu_2))^T(\mu_2 - (r_1\mu_1 + r_2\mu_2)) \\
&= (r_2\mu_1 - r_2\mu_2)^T(r_2\mu_1 - r_2\mu_2) + (r_1\mu_2 - r_1\mu_1)^T(r_1\mu_2 - r_1\mu_1) \\
&= r_2^2(\mu_1 - \mu_2)^T(\mu_1 - \mu_2) + r_1^2(\mu_2 - \mu_1)^T(\mu_2 - \mu_1) \\
&= (r_1^2 + r_2^2)(\mu_1 - \mu_2)^T(\mu_1 - \mu_2) \\
&\propto \Sigma_{betweenII}
\end{aligned}
$$

$$
where \qquad r_{1,2} = N_{1,2}/(N_1 + N_2)
$$

---

### Remark for PCA
Proof of $tr(DA) \leq tr(D)tr(A)$ given $D$ is diagonal.

$$
\begin{aligned}
(DA)_{nm} &= (D_{nm}A_{nm}) && \text{if D is doagonal} \\
tr(DA) &= \sum_{n=1}^{N}(DA)_{nn} \\
(tr(DA))^2 &= \left(\sum_{n=1}^{N}(DA)_{nn}\right)^2 \\
&\leq \left(\sum_{n=1}^{N}D_{nn}^2\right)\left(\sum_{n=1}^{N}A_{nn}^2\right) && \text{Cauchy Schwarz states } (\bar{x}\bar{y})^2 \leq |\bar{x}|^2|\bar{y}|^2 \text{ or equivalently } (\sum_n x_n y_n)^2 \leq (\sum_n x_n^2)(\sum_n y_n^2) \\
&\leq \left(\sum_{n=1}^{N}D_{nn}\right)^2\left(\sum_{n=1}^{N}A_{nn}\right)^2 && \text{sum of square is smaller than square of sum (this is proved by induction on n)} \\
&= (tr(D))^2(tr(A))^2 \\
tr(DA) &\leq tr(D)tr(A)
\end{aligned}
$$

### Remark for LDA
The inverse of a symmetric matrix is also symmetric, here is the proof :

$$
\begin{aligned}
AA^{-1} &= I \\
A^T A^{-1} &= I \\
(A^T A^{-1})^T &= I \\
A^{-T}A &= I \\
A^{-T} = (A^{-1})^T &= A^{-1} && \text{inverse of symmetric matrix is symmetric}
\end{aligned}
$$

However the product of two symmetric matrices, including $(\Sigma_{within})^{-1}\Sigma_{between}$, is not necessarily symmetric (find counter example).

## 3.4 Neural Network

A partite in a graph is a set of vertices having no edges connecting any two vertices in the set. Graphs which can be partitioned into multiple partites are called multipartite-graph. Network is a specialized multipartite-graph in which each partite is allowed to have edges connecting to at most two other neighbouring partites.

### 3.4a Topology of neural network `5 = [L,M0,M1,K,N]`

- •1   $L+1$ layers, indiced as $l \in [0,L]$
- •   layer $0$ is input layer, on which, each node is simply a scalar value $x$    } *about layer*
- •   layer $[1,L]$ are hidden/output layer, on which, each node is a complicated `struct` $\{ y, x, z, e, \delta \}$
- •23   there are $M_0$ input nodes, $M_1$ output nodes
- •4   there are $K$ nodes in whole network, nodes are indiced as $k \in [1,K]$   } *about node*
- •   nodes are not labelled with layer index
- •   there exist edges between layer $l$ and layer $l+1$ only, $\forall l \in [0,L-1]$
- •   $w_{ij}$ is weight of edge linking node $i$ in layer $l$ and node $j$ in layer $l+1$, for $\forall l \in [0,L-1]$   } *about edge*
- •   weights are tuned so as to minimise training error uising gradient descent
- •   key is to find gradient $\partial E/\partial w_k$ for all weights efficiently, without duplicated calculations
- •   network function is denoted as $B = f(A|w)$, where $A \in \Re^{M0}$ and $B \in \Re^{M1}$
- •5   network weights are trained with $N$ data points ($A_n \in \Re^{M0}, B_n \in \Re^{M1}$) $\forall n \in [1,N]$, error $E \in \Re^1$ is defined as :

$$E \quad = \quad \frac{1}{2}\sum_{n=1}^{N}\left\| f(A_n \mid w) - B_n \right\|^2 \qquad \textit{A and B are input/output of network fct}$$

$$= \quad \frac{1}{2}(F(A|w)-B)^T (F(A|w)-B) \qquad \textit{where } F(A|w) = \begin{bmatrix} f(A_1 \mid w) \\ f(A_2 \mid w) \\ ... \end{bmatrix}$$

### 3.4b Topology of a node   `5 = [y,x,z,e,δ]` and 5 steps

A node can be visualized with $x$ as node value and $w$ as edge value *(hide y, z, e and $\delta$ for clarity)*.



*For node k in layer l*

| | | | |
|---|---|---|---|
| 1 | *node value* | = | $x_k$ |
| 2 | *node delta (we will define it later)* | = | $\delta_k$ |
| 3 | *combo of node values from layer l-1* | = | $y_k$ |
| 4 | *combo of node deltas from layer l+1* | = | $e_k$ |
| | *s is activation fct, such that*   $s(x)$ | = | $1/(1+e^{-x})$ |
| 5 | *activation slope*    $z_k$ | = | $s'(y_k)$ |

**Step1** Initialize value in layer $0$ :

$$\{x_k\} \quad = \quad A_n \qquad \textit{where } \{x_k\} \textit{ is all } M_0 \textit{ node-values in layer 0}$$

**Step2** Forward propagation of input :

$$y_k \quad = \sum_{i \in layer[l-1]} w_{ik}x_i + \theta_k \qquad \textit{by definition}$$

$$x_k \quad = \quad s(y_k) \qquad \textit{by definition}$$

$$z_k \quad = \quad s'(y_k) \qquad \textit{by definition}$$

**Step3** Calculate error in layer $L$ :

$$\{\delta_k\} \quad = \quad f(A_n \mid w) - B_n \qquad \textit{where } \{\delta_k\} \textit{ is all } M_1 \textit{ node-errors in layer L}$$

**Step4** Backward propagation of error by chain rule :

$$e_k \quad = \sum_{j \in layer[l+1]} w_{kj}\delta_j \qquad \textit{by definition}$$

$$\delta_k \quad = \quad z_k e_k \qquad \textit{(equation 1) to be proved in next section}$$

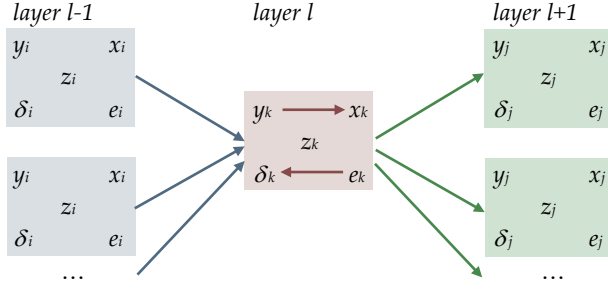$$\frac{\partial E}{\partial w_{ik}} \quad = \quad x_i \delta_k \qquad \textit{(equation 2) to be proved in next section}$$

We can do online training *(sample-by-sample weight update)* or offline training *(batch weight update)* :

$$w_{ik}' \quad = \quad w_{ik} - c\frac{\partial E}{\partial w_{ik}} \qquad\qquad \text{for all node } k \text{ in layer } l \text{ and all node } i \text{ in layer } l\text{-}1, \forall l \in [1,L]$$

$$= \quad w_{ik} - cx_i\delta_k$$

$$= \quad w_{ik} - cx_i z_k e_k$$

### 3.4c Proof of equation 1-2

The idea of backpropagation is to breakdown the gradient by chain rule into a two-step recursive function. Let's consider node $k$ in layer $l$, node $i$ in layer $l$-1 and node $j$ in layer $l$+1.



If we consider the network as a DAG, we can differentiate a variable v wrt another variable u, as long as there exists a path that starts from v and reaches u. Don't do it the opposite way.

Firstly, we define $\delta_k$ and apply chain rule to the subgraph started from node $k$ in layer $l$ :

$$\delta_k \quad \equiv \quad \frac{\partial E}{\partial y_k} \qquad\qquad \textit{objective is to expand } \frac{\partial output}{\partial input}$$

$$= \quad \frac{\partial E}{\partial x_k}\frac{\partial x_k}{\partial y_k} \qquad\qquad \textit{since a shock in } y_k \textit{ trigger a shock only in } x_k$$

$$= \quad \left[\sum_{j\in layer[l+1]}\frac{\partial E}{\partial y_j}\frac{\partial y_j}{\partial x_k}\right]\frac{\partial x_k}{\partial y_k} \qquad \textit{since a shock in } x_k \textit{ trigger a shock only in all } y_j \textit{ in layer } l+1$$

$$= \quad \left[\sum_{j\in layer[l+1]}\frac{\partial E}{\partial y_j}w_{kj}\right]z_k \qquad \textit{since } \frac{\partial y_j}{\partial x_k} = \frac{\partial \sum_{\eta \in layer[l]}w_{\eta j}x_\eta}{\partial x_k} = w_{kj} \textit{ and } \frac{\partial x_k}{\partial y_k} = z_k$$

$$= \quad [\sum_{j\in layer[l+1]}\delta_j w_{kj}]z_k \qquad \textit{since } \frac{\partial E}{\partial y_j} = \delta_j$$

$$= \quad z_k e_k \qquad\qquad \textit{(1) is proved}$$

Secondly, we solve the derivative of error wrt each weight and wrt each bias term :

$$\frac{\partial E}{\partial w_{ik}} \quad = \quad \frac{\partial E}{\partial y_k}\frac{\partial y_k}{\partial w_{ik}}$$

$$= \quad \delta_k \frac{\partial}{\partial w_{ik}}(\sum_{i\in layer[l-1]}w_{ik}x_i + \theta_k) \qquad \textit{since } \frac{\partial E}{\partial y_k} = \delta_k$$

$$= \quad x_i\delta_k \qquad\qquad \textit{(2) is proved for weight}$$

$$\frac{\partial E}{\partial \theta_k} \quad = \quad \frac{\partial E}{\partial y_k}\frac{\partial y_k}{\partial \theta_k}$$

$$= \quad \delta_k \frac{\partial}{\partial \theta_k}(\sum_{i\in layer[l-1]}w_{ik}x_i + \theta_k)$$

$$= \quad \delta_k \qquad\qquad \textit{(2) is proved for bias}$$

## 3.5 Support Vector Machine

Hyperplane in $\Re^M$ can be specified as $xw + b = 0$, where $w$ (column matrix) is the normal vector of hyperplane, and $x$ (row matrix) is any point on the hyperplane. It is **NOT** necessary that $w^T w = 1$, then we have :

$$
\begin{aligned}
xw &= -b \\
|x||w|\cos\theta &= -b \\
\underbrace{|x|\cos\theta}_{\text{projection of } x \text{ on } w} &= \underbrace{-b/|w|}_{\text{normal distance}}
\end{aligned}
$$

where $\theta$ is included angle between point $x$ and normal vector $w$

10 remarks for three next sections ...

```
(1) Formulation          (4) SVM objective    (8) How to solve b?
(2) Introduce Lagrange   (5) SVM constraint   (9) How to solve y?
(3) Differentiate wrt w&b (6) Definition of H  (0) Kernel trick
                         (7) Definition of α
```

### 3.5a Linear separable classification 10

Given $N$ data points $x_n$ $n \in [1,N]$ in $\Re^M$ and boolean labels $y_n \in [-1,+1]$, we find a double-hyperplane with maximum margin such that all data points are correctly classified. Double hyperplane in $\Re^M$ is defined by $M+1$ parameters, $w \in \Re^M$ and $b \in \Re^1$ :

$$
\begin{aligned}
xw + b &= \pm 1 \\
margin &= \frac{-b+1}{|w|} - \frac{-b-1}{|w|} = \frac{2}{|w|}
\end{aligned}
$$

we do **NOT** assume $w^T w = 1$

$$
\max_{w,b} \frac{2}{|w|}
$$

such that $\begin{bmatrix} x_n w + b \geq +1 & for & y_n = +1 \\ x_n w + b \leq -1 & for & y_n = -1 \end{bmatrix}$

$$
\Rightarrow \quad \min_{w,b} \frac{w^T w}{2}
$$

such that $\begin{bmatrix} y_n(x_n w + b) \geq 1 & for & y_n = +1 \\ y_n(x_n w + b) \geq 1 & for & y_n = -1 \end{bmatrix}$ (1)

$$
\Rightarrow \quad \min_{w,b} \max_{\alpha} \frac{w^T w}{2} - \sum_n \alpha_n [y_n(x_n w + b) - 1]
$$

such that $\alpha_n \geq 0, \forall n \in [1,N]$ (2)

$$
\Rightarrow \quad \min_{w,b} \max_{\alpha} \underbrace{\frac{w^T w}{2} - (\sum_n \alpha_n y_n x_n)w - (\sum_n \alpha_n y_n)b + \sum_n \alpha_n}_{L}
$$

such that $\alpha_n \geq 0, \forall n \in [1,N]$

We can differentiate above objective function with respect to $w$ $b$ and $\alpha$, therefore setting up $M+1+N$ equations of $M+1+N$ unknowns. Yet in most *SVM* derivations, we differentiate with respect to $w$ and $b$ only, as we want to get an objective purely in $\alpha$ (no $w$ nor $b$).

$$
\begin{aligned}
\frac{\partial L}{\partial w} &= w^T - \sum_n \alpha_n y_n x_n = 0 \quad \Rightarrow \quad \sum_n \alpha_n y_n x_n = w^T \\
\frac{\partial L}{\partial b} &= \sum_n \alpha_n y_n = 0 \quad \Rightarrow \quad \sum_n \alpha_n y_n = 0
\end{aligned}
$$

Substituting into $L$, we have (term highlighted in blue is removed by $\sum_n \alpha_n y_n = 0$) :

$$
\max_{\alpha} \frac{(\sum_n \alpha_n y_n x_n)(\sum_n \alpha_n y_n x_n)^T}{2} - (\sum_n \alpha_n y_n x_n)(\sum_n \alpha_n y_n x_n)^T + \sum_n \alpha_n
$$

such that $\alpha_n \geq 0, \forall n \in [1,N]$ and $\sum_n \alpha_n y_n = 0$ (3)

$$
\Rightarrow \quad \max_{\alpha} -\frac{1}{2}(\sum_n \alpha_n y_n x_n)(\sum_n \alpha_n y_n x_n)^T + \sum_n \alpha_n
$$

such that $\alpha_n \geq 0, \forall n \in [1,N]$ and $\sum_n \alpha_n y_n = 0$

$$
\Rightarrow \quad \max_{\alpha} -\frac{1}{2}\sum_{n,m} \alpha_n \alpha_m (y_n x_n)(y_m x_m)^T + \sum_n \alpha_n
$$

such that $\alpha_n \geq 0, \forall n \in [1,N]$ and $\sum_n \alpha_n y_n = 0$

$$
\Rightarrow \quad \max_{\alpha} -\frac{1}{2}\alpha^T H\alpha + l^T \alpha
$$

*Negative quadratic term coeff, hence maximization (as opposed to minimization)*

such that $\alpha_n \geq 0, \forall n \in [1,N]$ and $y^T \alpha = 0$ (4)&(5)

$$
\begin{aligned}
H &= [H_{nm}] \\
H_{nm} &= y_n y_m x_n x_m^T \\
\alpha &= \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ ... \end{bmatrix}
\end{aligned}
$$

$H$ is $N \times N$ matrix of dot products (6)

$\alpha$ is column matrix of Lagrange multiplers (7)

$y$ is column matrix of labels

$l$ is column matrix of all one

27

The algorithm can be summarised as :

- solve $\alpha$ by quadratic programming (such as *sequential quadratic programming* or *sequential minimal optimization*)
- solve $b$ by picking any support vector such that :

$$
\begin{aligned}
b &= y_{sv} - x_{sv}w \\
&= y_{sv} - x_{sv}(\textstyle\sum_n \alpha_n y_n x_n)^T \\
&= y_{sv} - \textstyle\sum_n \alpha_n y_n (x_{sv}x_n^T)
\end{aligned}
\tag{8}
$$

- there is no need to solve for $w$, in testing we simply apply :

$$
\begin{aligned}
y &= xw + b \\
&= x(\textstyle\sum_n \alpha_n y_n x_n)^T + b \\
&= \textstyle\sum_n \alpha_n y_n (xx_n^T) + b
\end{aligned}
\tag{9}
$$

- it can be extended to nonlinear hyperplane using the kernel trick, i.e. replacing all dot product $x_n x_m^T$ by kernel :

$$
\begin{aligned}
H_{nm} &= y_n y_m K(x_n, x_m) \qquad\qquad && where\ K(x_n, x_m) : (\Re^M, \Re^M) \to \Re^1 \\
b &= y_{sv} - \textstyle\sum_n \alpha_n y_n K(x_n, x_{sv}) \\
y &= \textstyle\sum_n \alpha_n y_n K(x_n, x) + b \\
&= \textstyle\sum_n \alpha_n y_n K(x_n, x) - \textstyle\sum_n \alpha_n y_n K(x_n, x_{sv}) + y_{sv} \\
&= \textstyle\sum_n \alpha_n y_n (K(x_n, x) - K(x_n, x_{sv})) + y_{sv}
\end{aligned}
\tag{10}
$$

*3.5b Linear non-separable classification*

By introducing an error term $\xi_n$ for each data point, and an error weight $C$, we have new objective *(new items are in red)* :

$$
\min_{w,b,\xi} \frac{w^T w}{2} + C\textstyle\sum_n \xi_n \qquad\qquad such\ that\ \begin{bmatrix} y_n(x_n w + b) \ge 1 - \xi_n & for & y_n = +1 \\ y_n(x_n w + b) \ge 1 - \xi_n & for & y_n = -1 \end{bmatrix} \tag{1}
$$

$$
and\ \xi_n \ge 0, \forall n \in [1, N]
$$

$$
\Rightarrow \quad \min_{w,b,\xi} \max_{\alpha,\beta} \frac{w^T w}{2} + C\textstyle\sum_n \xi_n - \textstyle\sum_n \alpha_n[y_n(x_n w + b) - 1 + \xi_n] - \textstyle\sum_n \beta_n \xi_n \qquad such\ that\ \alpha_n, \beta_n \ge 0, \forall n \in [1, N] \tag{2}
$$

$$
\Rightarrow \quad \min_{w,b} \max_{\alpha} \frac{w^T w}{2} + C\textstyle\sum_n \xi_n + \begin{bmatrix} -(\sum_n \alpha_n y_n x_n)w - (\sum_n \alpha_n y_n)b \\ + \sum_n \alpha_n - \sum_n \alpha_n \xi_n - \sum_n \beta_n \xi_n \end{bmatrix} \qquad such\ that\ \alpha_n, \beta_n \ge 0, \forall n \in [1, N]
$$

We can differentiate above objective function with respect to $w$, $b$ and $\xi$.

$$
\begin{aligned}
\frac{\partial L}{\partial w} &= w^T - \textstyle\sum_n \alpha_n y_n x_n &= 0 \quad &\Rightarrow \quad &\textstyle\sum_n \alpha_n y_n x_n &= w^T \\
\frac{\partial L}{\partial b} &= \textstyle\sum_n \alpha_n y_n &= 0 \quad &\Rightarrow \quad &\textstyle\sum_n \alpha_n y_n &= 0 \\
\frac{\partial L}{\partial \xi_n} &= C - \alpha_n - \beta_n &= 0 \quad &\Rightarrow \quad &\alpha_n + \beta_n &= C
\end{aligned}
$$

Substituting into $L$, we have the **same objective** as separable case, there are differences in constraints only. The terms highlighted in red are all cancelled by $\alpha_n + \beta_n = C$ .

$$
\max_{\alpha} \frac{(\sum_n \alpha_n y_n x_n)(\sum_n \alpha_n y_n x_n)^T}{2} - (\textstyle\sum_n \alpha_n y_n x_n)(\textstyle\sum_n \alpha_n y_n x_n)^T + \textstyle\sum_n \alpha_n \qquad such\ that\ 0 \le \alpha_n \le C, \forall n \in [1, N]\ and\ \textstyle\sum_n \alpha_n y_n = 0 \tag{3}
$$

$$
\Rightarrow \quad \dots
$$

$$
\Rightarrow \quad \max_{\alpha} -\frac{1}{2}\alpha^T H \alpha + l^T \alpha \qquad\qquad such\ that\ 0 \le \alpha_n \le C, \forall n \in [1, N]\ and\ y^T \alpha = 0 \tag{4\&5}
$$

where $H$, $\alpha$, $b$ and $y$ are defined and calculated in the same way as linear separable case. Kernel applies too. $\qquad$ (6)-(10)

## 3.5c Linear regression

By introducing an insenstive tube of $\pm\varepsilon$, error term $\xi_n$ is considered as regression error outside the tube.

$$\min_{w,b,\xi} \frac{w^T w}{2} + C\sum_n (\xi_n^+ + \xi_n^-) \qquad where \quad \begin{bmatrix} x_n w + b + \varepsilon \geq y_n - \xi_n^+ \\ x_n w + b - \varepsilon \leq y_n + \xi_n^- \end{bmatrix} \qquad such\ that\ \ \xi_n^\pm \geq 0, \forall n \in [1,N] \qquad (1)$$

$$\Rightarrow \quad \min_{w,b,\xi^\pm} \max_{\alpha^\pm, \beta^\pm} \frac{w^T w}{2} + C\sum_n (\xi_n^+ + \xi_n^-) - \begin{bmatrix} \sum_n \alpha_n^+[(x_n w + b) - y_n + \varepsilon + \xi_n^+] + \sum_n \beta_n^+ \xi_n^+ \\ \sum_n \alpha_n^-[y_n - (x_n w + b) + \varepsilon + \xi_n^-] + \sum_n \beta_n^- \xi_n^- \end{bmatrix} \qquad such\ that\ \ \alpha_n^\pm, \beta_n^\pm \geq 0, \forall n \in [1,N] \qquad (2)$$

$$\Rightarrow \quad \min_{w,b,\xi^\pm} \max_{\alpha^\pm, \beta^\pm} \frac{w^T w}{2} + C\sum_n (\xi_n^+ + \xi_n^-) - \begin{bmatrix} +(\sum_n \alpha_n^+ x_n)w + (\sum_n \alpha_n^+)b + \sum_n \alpha_n^+(-y_n + \varepsilon) + \sum_n (\alpha_n^+ + \beta_n^+)\xi_n^+ \\ -(\sum_n \alpha_n^- x_n)w - (\sum_n \alpha_n^-)b + \sum_n \alpha_n^-(+y_n + \varepsilon) + \sum_n (\alpha_n^- + \beta_n^-)\xi_n^- \end{bmatrix} \qquad such\ that\ \ \alpha_n^\pm, \beta_n^\pm \geq 0, \forall n \in [1,N] \qquad$$

We can differentiate above objective function with respect to $w$, $b$ and $\xi^\pm$.

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| $\frac{\partial L}{\partial w}$ | $=$ | $w^T - \sum_n (\alpha_n^+ - \alpha_n^-)x_n$ | $=$ | $0$ | $\Rightarrow$ | $\sum_n (\alpha_n^+ - \alpha_n^-)x_n$ | $= \quad w^T$ |
| $\frac{\partial L}{\partial b}$ | $=$ | $\sum_n (\alpha_n^+ - \alpha_n^-)$ | $=$ | $0$ | $\Rightarrow$ | $\sum_n (\alpha_n^+ - \alpha_n^-)$ | $= \quad 0$ |
| $\frac{\partial L}{\partial \xi_n^+}$ | $=$ | $C - (\alpha_n^+ + \beta_n^+)$ | $=$ | $0$ | $\Rightarrow$ | $\alpha_n^+ + \beta_n^+$ | $= \quad C$ |
| $\frac{\partial L}{\partial \xi_n^-}$ | $=$ | $C - (\alpha_n^- + \beta_n^-)$ | $=$ | $0$ | $\Rightarrow$ | $\alpha_n^- + \beta_n^-$ | $= \quad C$ |

Substituting into $L$, we have term in blue cancelled by $\sum_n (\alpha_n^+ - \alpha_n^-) = 0$ and terms in red cancelled by $\alpha_n^\pm + \beta_n^\pm = C$.

$$\min_w \max_{\alpha^\pm} \frac{w^T w}{2} - \begin{bmatrix} +(\sum_n \alpha_n^+ x_n)w + \sum_n \alpha_n^+(-y_n + \varepsilon) \\ -(\sum_n \alpha_n^- x_n)w + \sum_n \alpha_n^-(+y_n + \varepsilon) \end{bmatrix} \qquad such\ that\ \ 0 \leq \alpha_n^\pm \leq C, \forall n \in [1,N]\ \ and\ \ l^T \alpha = 0 \qquad (3)$$

$$\Rightarrow \quad \min_w \max_{\alpha^\pm} \frac{w^T w}{2} - \sum_n (\alpha_n^+ - \alpha_n^-)x_n w + \sum_n (\alpha_n^+ - \alpha_n^-)y_n - \sum_n (\alpha_n^+ + \alpha_n^-)\varepsilon \qquad such\ that\ \ 0 \leq \alpha_n^\pm \leq C, \forall n \in [1,N]\ \ and\ \ l^T \alpha = 0$$

$$\Rightarrow \quad \max_{\alpha^\pm} -\frac{(\sum_n (\alpha_n^+ - \alpha_n^-)x_n)(\sum_n (\alpha_n^+ - \alpha_n^-)x_n)^T}{2} + \sum_n (\alpha_n^+ - \alpha_n^-)y_n - \sum_n (\alpha_n^+ + \alpha_n^-)\varepsilon \qquad such\ that\ \ 0 \leq \alpha_n^\pm \leq C, \forall n \in [1,N]\ \ and\ \ l^T \alpha = 0$$

$$\Rightarrow \quad \max_{\alpha^\pm} -\frac{1}{2}\alpha^T H\alpha + y^T \alpha - l^T \alpha * \cdot \varepsilon \qquad such\ that\ \ 0 \leq \alpha_n^\pm \leq C, \forall n \in [1,N]\ \ and\ \ l^T \alpha = 0 \qquad (4)\&(5)$$

| | | | |
|---|---|---|---|
| $where$ | $H$ | $=$ | $[H_{nm}]$ |
| | $H_{nm}$ | $=$ | $x_n x_m^T$ | (6) |

$$\alpha = \begin{bmatrix} \alpha_1^+ - \alpha_1^- \\ \alpha_2^+ - \alpha_2^- \\ \alpha_3^+ - \alpha_3^- \\ ... \\ \alpha_N^+ - \alpha_N^- \end{bmatrix} \qquad and \qquad \alpha* = \begin{bmatrix} \alpha_1^+ + \alpha_1^- \\ \alpha_2^+ + \alpha_2^- \\ \alpha_3^+ + \alpha_3^- \\ ... \\ \alpha_N^+ + \alpha_N^- \end{bmatrix} \qquad (7)$$

Finally solve $b$ by picking any support vector such that :

$$b = \begin{bmatrix} y_{sv} - x_{sv}w - \varepsilon - \xi_{sv}^+ \\ y_{sv} - x_{sv}w + \varepsilon + \xi_{sv}^- \end{bmatrix}$$

$$= \begin{bmatrix} y_{sv} - x_{sv}(\sum_n (\alpha_n^+ - \alpha_n^-)x_n)^T - \varepsilon - \xi_{sv}^+ \\ y_{sv} - x_{sv}(\sum_n (\alpha_n^+ - \alpha_n^-)x_n)^T + \varepsilon + \xi_{sv}^- \end{bmatrix}$$

$$= \begin{bmatrix} y_{sv} - \sum_n (\alpha_n^+ - \alpha_n^-)(x_{sv}x_n^T) - \varepsilon - \xi_{sv}^+ \\ y_{sv} - \sum_n (\alpha_n^+ - \alpha_n^-)(x_{sv}x_n^T) + \varepsilon + \xi_{sv}^- \end{bmatrix} \qquad (8)$$

| | | |
|---|---|---|
| $y$ | $=$ | $xw + b$ |
| | $=$ | $x(\sum_n (\alpha_n^+ - \alpha_n^-)x_n)^T + b$ |
| | $=$ | $\sum_n (\alpha_n^+ - \alpha_n^-)(xx_n^T) + b$ | (9) |

*Please read "Support Vector Machines Explained" by Tristan Fletcher*

## 3.6 Generative Adversarial Net

Given a set of training data $x_n \in \Re^M \ \forall n \in [1,N]$ drawn from distribution $p_{data}(x)$, *GAN* is a learning algorithm, which approximates the distribution $p_{data}(x)$ with a pair of neutral networks ($G$ and $D$) combating each other in a multistage game. $G$ is a ***generative network*** which generates output $x = G(z|w_G)$ with a random number generator $z$ as input. $D$ is a ***discriminative network*** which estimates the probability $y = D(x|w_D) \in [0,1]$ that an input $x$ is a sample drawn from $p_{data}(x)$ rather than generated by $G(z|w_G)$. In short :

- $G$ is a counterfeit generator, which performs exact imitation
- $D$ is a counterfeit detector, which differentiates training data *(hopefully y ~ 1)* from counterfeit data *(hopefully y ~ 0)*
- $G$ has $L$ input nodes, for $z \in \Re^L$ is a $L$ dimensional random number generator
- $G$ has $M$ output nodes, for $x \in \Re^M$ is a $M$ dimensional simulated data
- $D$ has $M$ input nodes, for $x \in \Re^M$ is a $M$ dimensional input data
- $D$ has 1 output node, for $y \in [0,1]$ is the probability that input data is drawn from $p_{data}(x)$
- weight sets for the two networks are $w_G$ and $w_D$ respectively
- distribution of random number generator $z$ is $p_{rng}(z)$
- distribution of counterfeit generator $G(z|w_G)$ is $p_G(x)$     *3 different pdfs*

The 3 distributions are related by :

$$p_{rng}(z)dz \quad = \quad p_G(x)dx \qquad\qquad \text{(1) suppose both L,M = 1}$$

$$= \quad p_G(x)\partial_z G(z|w_G)dz$$

$$p_G(x) \quad = \quad \frac{p_{rng}(z)}{\partial_z G(z|w_G)}$$

*and we hope* $\quad p_G(x) \quad \sim \quad p_{data}(x)$

*adjust generator's weight and RNG z so as to...*

*achieve exact imitation of training data*

Furthermore from *(1)*, we have :

$$f(x)p_G(x)dx \quad = \quad f(x)p_{rng}(z)dz \qquad\qquad \text{\textit{for any function f(x)}}$$

$$= \quad f(G(z|w_G))p_{rng}(z)dz \qquad \text{\textit{since}} \ x = G(z|w_G)$$

$$\int f(x)p_G(x)dx \quad = \quad \int f(G(z|w_G))p_{rng}(z)dz$$

$$E_{x \sim p_G}[f(x)] \quad = \quad E_{x \sim p_{rng}}[f(G(z|w_G))] \qquad\qquad \text{(2)}$$

## GAN concept

Log likelihood of correct discrimination is :

$$\ln(D(x|w_D)) \qquad\qquad \text{\textit{for all x drawn from }} p_{data}(x)$$

$$\ln(1 - D(x|w_D)) \qquad\qquad \text{\textit{for all x drawn from }} p_G(x)$$

Discriminator maximizes the likelihood *wrt $w_D$*, while generator minimizes ***discriminator's optimal*** *wrt $w_G$*.

$$\min_{w_G}\max_{w_D} E_{x \sim p_{data}}[\ln(D(x|w_D))] + E_{x \sim p_G}[\ln(1 - D(x|w_D))]$$

$$= \quad \min_{w_G}\max_{w_D} E_{x \sim p_{data}}[\ln(D(x|w_D))] + E_{z \sim p_{rng}}[\ln(1 - D(G(z|w_G)|w_D))] \qquad \text{\textit{applying (2)}}$$

$$= \quad \min_{w_G}\max_{w_D}\left( \int_\infty p_{data}(x)\ln(D(x|w_D))dx + \int_\infty p_{rng}(z)\ln(1 - D(G(z|w_G)|w_D))dz \right) \qquad \text{\textit{minmax f(x,y)} $\neq$ \textit{maxmin f(x,y) consider matrix example}}$$
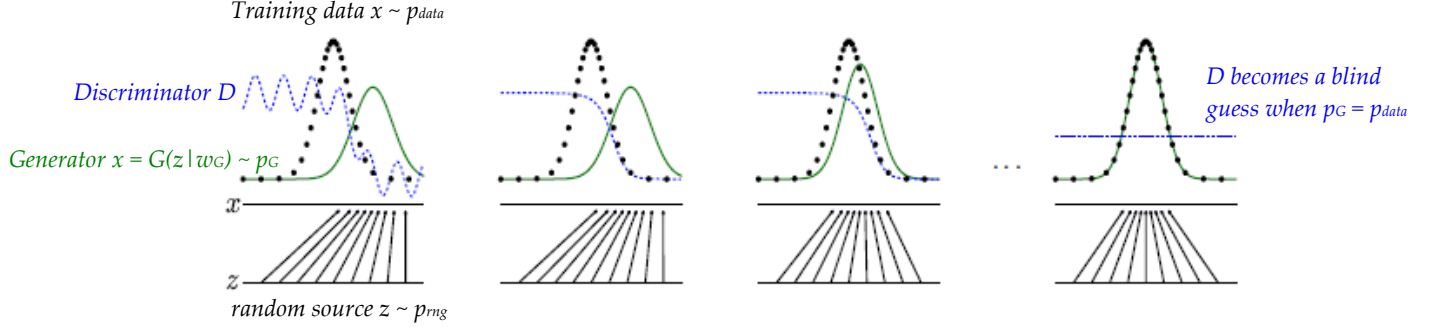
## GAN implementation

GAN is implemented as *K discriminator-iterations* plus *1 generator-iteration*.

```
while(!converge)
{    for(int k=0; k!=K ++k)                            different gradients
    {
        draw N samples from Pdata(x)
        draw N samples from Prng(z) and plug in G(z|wG)
        perform (stochastic) gradient ascent of ∇_wD ∑ln(D(x|wD)) + ∇_wD ∑ln(1 - D(G(z|wG)|wD))     [D-step]
    }
    draw N samples from Prng(z) and plug in G(z|wG)
    perform (stochastic) gradient descent of ∇_wG ∑ln(1 - D(G(z|wG)|wD))                              [G-step]
}
```

## GAN convergence

The algorithm converges as illustrated, the approximation by $G$ gets improved, discrimination by $D$ becomes a blind guess.



*Training data $x \sim p_{data}$*

*Discriminator D*

*Generator $x = G(z \mid w_G) \sim p_G$*

$x$

$z$

*random source $z \sim p_{rng}$*

*D becomes a blind guess when $p_G = p_{data}$*

**[D-step]** Given $G(z \mid w_G)$ in previous iteration, optimal of $D$ in this iteration is given by :

$$\underset{w_D}{\arg\max} \int_\infty [p_{data}(x)\ln(D(x \mid w_D)) + p_G(x)\ln(1 - D(x \mid w_D))]dx$$

$\rightarrow \quad \int_\infty [p_{data}(x)\nabla_{w_D}\ln(D(x \mid w_D)) + p_G(x)\nabla_{w_D}\ln(1 - D(x \mid w_D))]dx \quad = \quad 0 \qquad$ *since $p_{data}(x)$ and $p_G(x)$ are independent on $w_D$*

$\rightarrow \quad p_{data}(x)\nabla_{w_D}\ln(D(x \mid w_D)) + p_G(x)\nabla_{w_D}\ln(1 - D(x \mid w_D)) \quad = \quad 0 \qquad$ *for all x, given enough degree of freedom in $w_D$*

$\rightarrow \quad D(x \mid w_D) \quad = \quad \dfrac{p_{data}(x)}{p_{data}(x) + p_G(x)} \qquad\qquad$ *see remark*

---

*Remark : find* $\underset{D}{\max}[a\ln D + b\ln(1 - D)]$

$\qquad\qquad 0 \quad = \quad \partial_D(a\ln D + b\ln(1 - D))$

$\qquad\qquad 0 \quad = \quad a/D + b/(1 - D)$

$\qquad\qquad bD \quad = \quad a(1 - D)$

$\qquad\qquad D \quad = \quad a/(a + b)$

---

**[G-step]** Given optimum $D$ in previous iteration, optimal of $G$ in this iteration is given by :

$$\underset{w_G}{\min}\left( \underset{w_D}{\max} E_{x \sim p_{data}}[\ln(D(x \mid w_D))] + E_{x \sim p_G}[\ln(1 - D(x \mid w_D))] \right)$$

$= \quad \underset{w_G}{\min}\left( \underset{w_D}{\max} E_{x \sim p_{data}}\left[ \ln\dfrac{p_{data}(x)}{p_{data}(x) + p_G(x)} \right] + E_{x \sim p_G}\left[ \ln\dfrac{p_G(x)}{p_{data}(x) + p_G(x)} \right] \right) \qquad$ *plug in result from D-step*

$= \quad \underset{w_G}{\min}\left( \underset{w_D}{\max} E_{x \sim p_{data}}\left[ \ln\dfrac{p_{data}(x)}{p_{data}(x) + p_G(x)} + \ln 2 \right] + E_{x \sim p_G}\left[ \ln\dfrac{p_G(x)}{p_{data}(x) + p_G(x)} + \ln 2 \right] \right) - \ln 4$

$= \quad \underset{w_G}{\min}\left( \underset{w_D}{\max} E_{x \sim p_{data}}\left[ \ln\dfrac{p_{data}(x)}{(p_{data}(x) + p_G(x))/2} \right] + E_{x \sim p_G}\left[ \ln\dfrac{p_G(x)}{(p_{data}(x) + p_G(x))/2} \right] \right) - \ln 4 \qquad$ *read KL div doc*

$= \quad \underset{w_G}{\min}\left( D_{KL}(p_{data}(x) \| \dfrac{p_{data}(x) + p_G(x)}{2}) + D_{KL}(p_G(x) \| \dfrac{p_{data}(x) + p_G(x)}{2}) \right) - \ln 4 \qquad$ *since $D_{KL}(P \| Q) = E_P\left[ \log_2 \dfrac{P}{Q} \right]$*

$= \quad -\ln 4 \qquad\qquad$ *since KL divergence must be non-negative*

Divergence is non-negative, thus minimum of above objective happens when :

$\qquad p_{data}(x) \quad = \quad \dfrac{p_{data}(x) + p_G(x)}{2}$

*i.e.* $\quad p_{data}(x) \quad = \quad p_G(x)$

In short, when *GAN* is done, generator distribution converges to training data distribution. Can we introduce other magic numbers such as 3&6 instead of 2&4? Of course not, we need to construct another measure to generate the KL divergence.

- *GAN* is liked *EM* as a bistep recursion :
  - *E-step*  is analogous to  *D-step*
  - *M-step*  is analogous to  *G-step*


- *GAN* is liked KL divergence :
  - maximization of expected log likelihood


- *GAN* is liked reversed Monte Carlo simulation :
  - *MC* finds function expectation of *F(z)* from samples of random variable *z*, given $p_Z(z)$
  - *GAN* finds parameters of *x = F(z|w)* from samples of random variable $x_{data}$, given $p_Z(z)$