

11

SMAI

Book: (Pattern Classification
↳ Puda, Hart & Stork)

- 1) Frequentist Approach → Assume there are a fixed set of parameters
- 2) Bayesian Approach → Varying parameters

* Classification & Clustering

- Preprocessing → Feature Extraction → Classification → Post Processing (Risk + Context)
 - (Invariance to translation, rotation & scaling)

Note: Training & Test data must be drawn from the same distribution

Classification

- Model, Decision boundary, cost.
 - Adjust this boundary.
 - Minimize this (Decision theory) cost.
- Generalization vs Overfitting.
- Post processing after classification can correct errors based on context.

Learning:

- 1) Supervised → Training data
 - 2) Unsupervised → ex. clustering
 - 3) Reinforcement → ex. RNNs
- (Use same data multiple times)

2/1

Probability Density Function:

- $P(a < x < b) = \int_a^b p(x) dx$
 \uparrow
Continuous probability function

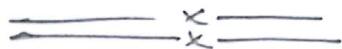
- Different types of functions.

$$\Rightarrow p(x) \geq 0$$

$$\Rightarrow \int_{-\infty}^{\infty} p(x) dx = 1$$

- Gaussian,

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



* Density Estimation:

- $P = \int_R p(x) dx$ } Probability of belonging to a region R .

\Rightarrow If R is small, so $p(x)$ is almost constant

$$P = \int_R p(x) dx \approx p(x) \int_R dx = p(x) \cdot V \rightarrow \text{"Volume" of region } R.$$

- So in 2D $\rightarrow R \approx (x_1, x_2) \rightarrow 2 \text{ Dimensional}$
we consider a small square/circle
so the area ("volume") $\rightarrow 0$.

\Rightarrow Consider we drew n data samples, $\{x_1, x_2, \dots\}$
if there are k points out of n lying in
a region R . ($P = k/n$ in that region)

* \Rightarrow But $P = p(x) \cdot V$

$$\therefore p(x) = \frac{k/n}{V}$$

} x is the point around
which we consider
region R .

- Since $V \rightarrow 0$, we would try make $n \rightarrow \infty$,
assume we have ∞ samples.



* Parzen Window → (Density Estimation)

1). Consider R to be a hypercube centered at \bar{x} .

- $h = \text{Length of edge}$

$$\Rightarrow \text{Volume of hypercube } (V) = h^d \xrightarrow{\text{Dimension of data}}$$

- Window function,

$$\phi\left(\frac{x_i - \bar{x}}{h}\right) = \begin{cases} 1, & \frac{|x_i - \bar{x}|}{h} \leq \frac{1}{2}, k=1,2,\dots,d \\ 0, & \text{otherwise} \end{cases}$$

\Rightarrow So total number of samples within R is k

$$k = \sum_{i=1}^n \phi\left(\frac{x_i - \bar{x}}{h}\right)$$

$$\therefore \text{Parzen prob density} \Rightarrow p(\bar{x}) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h^d} \phi\left(\frac{x_i - \bar{x}}{h}\right)$$

Note:

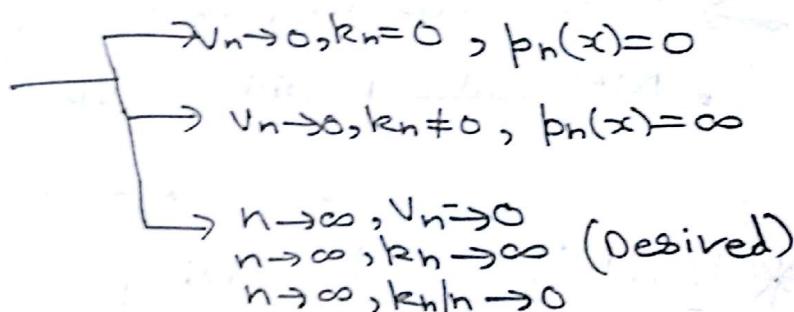
\Rightarrow We can use other shapes like spheres etc instead of hypercubes.

Parzen Window Dichotomizer:

- Classify into 2 classes $\rightarrow +1 \& -1$.
 \Rightarrow Consider parzen window $\xrightarrow{\text{new point } \bar{x}}$ if majority of neighbours are +1, then this point goes to +1 etc.

Convergence of Parzen Window:

$$\bullet p_n(\bar{x}) = \frac{k_n}{h \cdot V_n}$$



\Rightarrow So we can reduce our window size as $n \uparrow$.

(kn = Varying, V = Const)

* Note : Parzen estimation \rightarrow Reduce volume of window based on h

KNN estimation \rightarrow Reduce volume of window by fixing number of samples enclosed by it
(kn = Const, V = Varying)

$\overbrace{\quad \quad \quad}^x$

* KNN \rightarrow (Density Estimation)

2)

- Consider classes $\rightarrow w_i$, total = c classes

$$\Rightarrow p(x, w_i) = \frac{(k_i/h)}{V} \quad \left. \begin{array}{l} k_i = \text{Samples in window} \\ \text{belonging to class} \\ i. \end{array} \right\}$$

$$\Rightarrow p(w_i | x) = \frac{p(x, w_i)}{\sum_{j=1}^c p(x, w_j)} \Rightarrow \frac{k_i}{k}$$

(k = Samples in region with volume V)

\Rightarrow We try fix $k_n = \beta / \delta^n$, For various values of β .

$\overbrace{\quad \quad \quad}^x$

* Nearest Neighbours Classifier (NN)

- Eager Learning \rightarrow Pre compute based on training data
 - Lazy Learning \rightarrow Given test, calculate then and there
- \rightarrow NN-Classifier is Lazy Learning
- NN rules partitions space into Voronoi cells
- $\overbrace{\quad \quad \quad}^x$

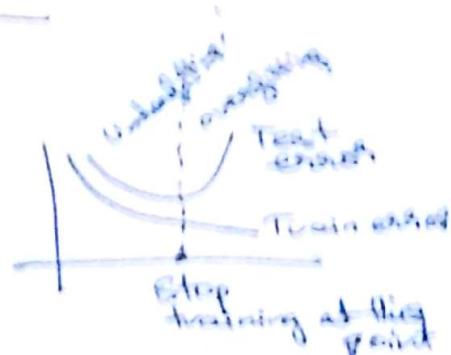
* k-NN Classifier

- Consider k nearest neighbours to classify.
- \rightarrow k should be odd to avoid ties.

- We could also use distance based weighting etc, so the target function is continuous.
- As the feature dimension increases, accuracy of the classifier drops.

* Classifier Evaluation:

Note:



1) Cross Validation

→ Random Subsampling:

- Consider some data as training & other as test

→ k -Fold Cross Validation:

- Consider windows of k samples as test data, and rest as training data
- $\Rightarrow n/k$ tests in total.

→ Leave 1 out cross validation:

- Use 1 sample for test
(In case our data is very small)

Note: Iris dataset → KNN Classifier,

Matlab/C++

K-fold cross validation

Tutorials: Wed 8:00-10:00

3

$$\bullet \text{Plane} \rightarrow (n - g_n) \cdot N = 0$$

\uparrow \uparrow \uparrow
 2 points on plane Normal

* Linear Discriminant Functions:

- We try find parameters of a LDF which minimizes the training error.

$$\Rightarrow g(x) = \begin{cases} > 0, \text{Class A} \\ < 0, \text{Class B} \\ = 0, \text{Boundary} \end{cases}$$

$$\Rightarrow g(x) = w^T x + w_0$$

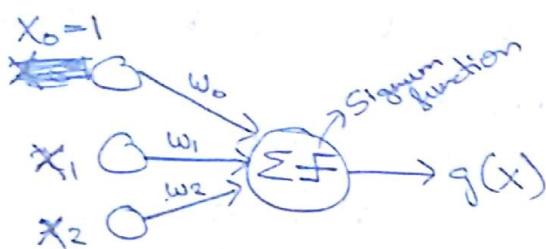
↳ Hyperplane to separate the classes

Does this actually represent a hyperplane
→ Proof below

→ In 2D, this becomes a line, $g(x) = mx + c$

— x —

* Perception:



$$w = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}$$

$$x = \begin{bmatrix} x_0 \\ \vdots \\ x_d \end{bmatrix}$$

$$\bullet g(x) = w_0 + w_1^T x_1 + w_2^T x_2 \quad \left. \begin{array}{l} \text{we will ignore the transpose and assume it exists from now on} \\ \text{This is the bias for the node, so we use } x_0 = 1 \end{array} \right\}$$

- We use the signum function so that if $g(x) \geq 0$, ~~$f_g(x)$~~ $f_g(x) = 1$, else $f_g(x) = 0$. (So we can split it into class A or class B)

Note:

$$\bullet g(x) = w^T x + w_0, \text{ consider this, } w = \text{Normal of hyperplane}$$

Consider an x ,

$$\Rightarrow x = x_p + \frac{w}{\|w\|} \rightarrow \text{Consider } x_p \text{ to be projection of } x \text{ on plane. } \frac{w}{\|w\|} = \text{Dist of } x \text{ from } x_p.$$

(Representation diagram in slides)

— x —

$$\begin{aligned} \cdot g(x) &= w^T(x_p + \frac{w}{\|w\|}) + w_0 \\ g(x) &= \underbrace{w^T x_p + w_0}_{\text{Since } x_p \text{ lies on plane}} + \frac{w^T w}{\|w\|} \end{aligned}$$

$$\therefore g(x) = g/\|w\|$$

$$\Rightarrow g_1 = \frac{g(x)}{\|w\|}, g_2 \propto g(x), g_0 = \frac{w_0}{\|w\|}$$

- So $g(x)$ can be used to show which side of the plane the point lies on, so which class it belongs to.

* \therefore It's a linear discriminant function its valid & represents a hyper plane.

\Rightarrow Generalized LDFs

$$\cdot g(x) = w^T x + w_0$$

\Rightarrow We can use this kind of $g(x)$ as an LDF since $g \propto g(x)$, as we proved above.

$$\Rightarrow g(x) = \sum_{i=0}^d w_i x_i, x_0 = 1 \quad \text{For the bias term.}$$

$$\Rightarrow g(x) = a^T Y, Y = \begin{bmatrix} x_0 \\ \vdots \\ x_d \end{bmatrix} = \begin{bmatrix} x_0 \\ X \end{bmatrix}$$

In this space, there is no bias, so this hyperplane goes through origin.

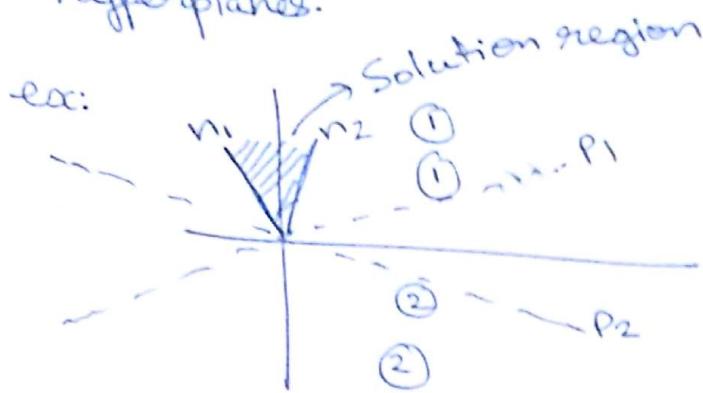
$$\Rightarrow Y = \phi(x)$$

- ϕ can be any arbitrary mapping which projects original data points to a higher dimension, hoping that we might find a linear hyperplane to divide the data. ex. $\begin{pmatrix} -2 \mapsto +1 \\ -1 \mapsto -1 \\ 2 \mapsto +1 \\ 1 \mapsto -1 \end{pmatrix}$
- Take square & we can divide \Rightarrow

Note: Curse of dimensionality, using \rightarrow higher dim spaces requires a lot more time on computation

Solution region: \rightarrow Slides to transform later on

- All possible solution normals for acceptable hyperplanes.



(Two category linear separability)

- We then take the dual form where a^T form the axes.

4

* Gradient Descent: (For linearly separable data)

- η = Learning rate. } the scale factor
- Scalar function $J(a)$ = Classification error described by a .
 \Rightarrow Minimize $J(a)$ using gradient descent.

*
$$a(k+1) = a(k) - \eta(k) \nabla J(a(k))$$

- Set a ϕ as a threshold, keep following the algorithm till $|J(a(k))| < \phi$, then stop.
or we could limit num of iterations, till $k \leq k_{\max}$

* Perceptron Criterion function:

- J should be a continuous function & not discrete.

$$\boxed{\text{J}_p(a) = \sum_{y \in Y} -a^T y}, \quad Y = \text{set of misclassified samples.}$$

(Piecewise linear)

($a^T y$ will be ≤ 0 for incorrect) $\rightarrow g(x) = a^T y$
 (Normalized inputs) $\qquad\qquad\qquad g(x) \leq 0$ for incorrect

• J_p is our criterion function.

$$\Rightarrow \nabla J_p(a) = \sum_{y \in Y} (-y)$$

$$\therefore \boxed{a(k+1) = a(k) + \eta(k) \sum_{y \in Y} y}$$

Set of
misclassified
points.

* This will only work (converge) if data is linearly separable.

\Rightarrow This is a Batch Perceptron.

Note: Single Sample Perceptron, iterate over each point doing, $a = a + \eta y_i$, if $a^T y_i < 0$ else we ignore it.

\Rightarrow Similar to batch except that we go one point at a time.

* Relaxation Procedures:

• Class of criterion functions & minimization methods

$$1) \Rightarrow \boxed{J_q(a) = \sum_{y \in Y} (a^T y)^2} \rightarrow \begin{array}{l} \text{Convergence to boundary} \\ \text{Dominated by longest sample.} \end{array}$$

(Squared error)

$$2) \boxed{J_m(a) = \frac{1}{2} \sum_{i \in X} \frac{(a^T y_i - b)^2}{m}} \rightarrow \begin{array}{l} \text{To stay} \\ \text{a little away} \\ \text{from the boundary} \end{array}$$

Minimum Squared Error: (Mainly for non linearly separable data)

- MSE considers all samples instead of just misclassified ones.

$$Y = [y_1, \dots, y_n]^T \quad \begin{matrix} \text{Normalized} \\ \text{inputs.} \end{matrix} \quad (\text{Parameters})$$

$$y_i = [y_{i0}, \dots, y_{id}]^T \quad , \quad a = [a_0, \dots, a_d]^T$$

We want,

$$\boxed{Ya = b}$$

$$(b > 0)$$

, $a = Y^{-1}b$ not possible,

since Y is rectangular & might be

The
dimensions
of Y
we want
from "soft"
hyperplane
hypothesis
are

Singular.

$$\Rightarrow \text{Error } e = \boxed{Ya - b} \quad \begin{matrix} \text{Normalized} \\ \text{error} \end{matrix}$$

$$\bullet \boxed{J(a) = \|e\|^2 = \|Ya - b\|^2} \quad \begin{matrix} \text{MSE} \\ \text{error} \\ \text{function.} \end{matrix}$$

Square of norm

$$\bullet \nabla J(a) = 2Y^T(Ya - b) \rightarrow \text{Set this to 0 for finding pseudo inverse}$$

Now,

$$\bullet Y^T Ya = Y^T b,$$

$$\therefore \boxed{a = (Y^T Y)^{-1} Y^T b} = \underbrace{Y^T b}_{\substack{\text{Pseudo inverse} \\ \text{of } Y.}}$$

When we achieve
0 error, so
 $\nabla J(a) = 0$ when
we have no
error.
(Minima)

$$\Rightarrow Y^T Y = \text{Square matrix \& often non-singular.}$$

If $Ya > 0$, i.e. $\forall i (a^T y_i) > 0$, then a hyperplane is always guaranteed. (For normalized y)

$$\Rightarrow Ya = b + \varepsilon, \text{ in reality,}$$

so some entries of b can be -ve.

$$|b_i| < |\varepsilon| \quad (\text{But we try to keep } b = 0)$$

∴ Even for linearly separable, we might not find a "perfect" hyperplane, but a reasonable

* Note: If no info is given, assume $b = [1 \dots 1]^T$

* When we take $ay > 0$ for the samples and $ay < 0$ for -ve samples, if we normalize then $ay > 0$ for all inputs. ($y_{-ve} \rightarrow -y_{-ve}$)

$\therefore ay_{-ve} \Rightarrow a(-y_{-ve}) > 0 \}$ for -ve inputs.

We can solve this easily if we just assume $ay = b$ since this is just an equation & not an inequality.

$$\therefore ya = b$$

$$(b > 0)$$

Normalized inputs.

Now if we manage to find an a which can do this then we are done

Using pseudo inverse, we can estimate a value for a . For estimated a^* , If $ya^* > 0$, then its a separating plane

* This method is a lot simpler and faster than gradient descent also works for non-linearly separable case.

Note: Some kinds of variations in b (Non uniform scaling of b) tonight help at times.
There are ways in which this is done.

Note: Normalized Inputs:

- * If $y \in Y$ is a +ve sample we use it as it is
- * If $y \in Y$ is a -ve sample we use $-y$.

$$\therefore ay_{-ve} \leq 0$$

$$a(-y_{-ve}) > 0$$

To always have $ay > 0$ we consider normalized inputs

$$*\text{ Widrow-Hoff } \quad a(k+1) = a(k) - \eta(k) \cdot Y^T(k) y(k)$$

Note:

Iterative Gradient Descent could be used with the MSE error function.

* Widrow-Hoff | LMS Procedure: (Stochastic)

- Using all our datapoints ^{at once} may be computationally very expensive.

\Rightarrow So we do,

$$a(k+1) = a(k) = -\eta(i) y_i (y_i^T a(k)) + \text{for } i = \text{mod}(i+1, n)$$

- Annealing $\eta(k) = \eta(i)/k$ for learning rate

\Rightarrow This is using GD along with MSE error on one sample at a time.

* Ho-Kashyap Procedure: (Vary b as well instead of fixing)

- Using modified gradient descent, minimize $J(a, b)$

- Dont fix $b = [1, \dots, 1]^T$, we use the MSE error here as well.

$$\bullet J(a, b) = \|e\|^2 = \|Y_a - b\|^2$$

\Rightarrow Fix b and minimize $J(a, b)$ wrt a .

* \Rightarrow Then fix a and minimize $J(a, b)$ wrt b .

- $b(k+1) = b(k) - \eta(k) \nabla_b J$
- $a(k+1) = a(k) - \eta(k) Y^T (Y_a(k) - b(k+1))$

$$\Rightarrow \nabla_b J = 2(Y_a - b)(-1) = -2e$$

$$\therefore e = -\frac{1}{2} \nabla_b J$$

- $\nabla_b J$ should always be -ve, so that we always keep adding to b , ie $b(k+1) \geq b(k)$
- \Rightarrow In this way we can ensure b is never going to be -ve.

∴ To do this we cant ensure $\nabla_b J < 0$,
so we do $\rightarrow \frac{1}{2}(\nabla_b J - |\nabla_b J|)$ instead of $\nabla_b J$

$$\begin{aligned}\therefore b(k+1) &= b(k) - \eta(k) \cdot \frac{1}{2}(\nabla_b J - |\nabla_b J|) \\ &= b(k) + \eta(k) \left(-\frac{1}{2} \nabla_b J + \frac{1}{2} |\nabla_b J| \right)\end{aligned}$$

$$b(k+1) = b(k) + \eta(k)(e + |e|)$$

- So b keeps increasing or remains equal.

\Rightarrow Pseudocode in slides.

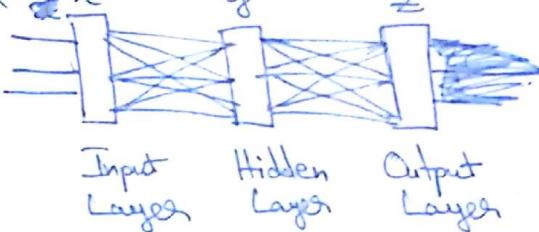
- Instead of updating $a(k+1)$ iteratively we could also just do $a(k+1) = Y^+ b(k+1)$
(Using pseudo inverse)

;

Neural Nets:

- Simple perceptrons can only handle linearly separable cases well. So extending to neural nets which can identify complex boundaries was done.
- Example of how non-linearity is modelled in slides.

* Backpropagation: (Derivation)



$$\bullet y_j = f(\text{net}_j) = \text{sgn}(\text{net}_j)$$

$$\bullet z_k = f(\text{net}_k) = \text{sgn}(\text{net}_k)$$

$$\bullet g(x) = z_k = \sum_{j=1}^{n_H} w_{kj} f\left(\sum_{i=1}^d x_i w_{ji} + w_{j0}\right) + w_{k0}$$

$$\Rightarrow J(w) = \frac{1}{2} \sum_{k=1}^c (t_k - z_k)^2 = \frac{1}{2} \|t - z\|^2$$

$$\Delta w = -\eta \frac{\partial J}{\partial w}, \Delta w_{pq} = -\eta \frac{\partial J}{\partial w_{pq}}$$

$$w(m+1) = w(m) + \Delta w$$

$$\frac{\partial J}{\partial w_{kj}} =$$

$$\frac{\partial J}{\partial w_{kj}} = \frac{\partial J}{\partial \text{net}_k} \cdot \frac{\partial \text{net}_k}{\partial w_{kj}} = \tilde{f}'(\delta_k) \cdot \frac{\partial \text{net}_k}{\partial w_{kj}} = \tilde{f}'(\delta_k) \cdot (-1) \cdot (t_k - z_k) \cdot f'(\text{net}_k) \cdot 1$$

$$-\delta_k = +\frac{\partial J}{\partial \text{net}_k} = +\frac{\partial J}{\partial z_k} \cdot \frac{\partial z_k}{\partial \text{net}_k}$$

$$= +1 \cdot \frac{1}{2} (t_k - z_k) (-1) \cdot f'(\text{net}_k)$$

$$-\delta_k = \cancel{+1} \cdot (t_k - z_k) \cdot f'(\text{net}_k)$$

$$\bullet \text{net}_j = \sum_{i=1}^d x_i w_{ji} + w_{j0}$$

$$\Rightarrow w_{ab} = \text{weight from } a \text{ to } b$$

$$\text{net}_j = \sum_{i=0}^d x_i w_{ji} = w_j^T x$$

$$\text{net}_k = \sum_{j=1}^{n_H} y_j w_{kj} + w_{k0}$$

$$\text{net}_k = \cancel{\sum_{j=0}^{n_H} y_j w_{kj}} \cancel{+ w_{k0}} = \sum_{j=0}^{n_H} y_j w_{kj} = w_k^T x$$

$$\cdot \frac{\partial \text{net}_k}{\partial w_{kj}} = y_j$$

$$\begin{aligned} \Rightarrow \frac{\partial J}{\partial w_{ji}} &= \frac{\partial J}{\partial \text{net}_j} \cdot \frac{\partial \text{net}_j}{\partial w_{ji}} \\ &= \frac{\partial J}{\partial y_j} \cdot \frac{\partial y_j}{\partial \text{net}_j} \cdot x_i \\ &= \frac{\partial J}{\partial y_j} \cdot f'(\text{net}_j) \cdot x_i \rightarrow ① \end{aligned}$$

Now,

$$\begin{aligned} \frac{\partial J}{\partial y_j} &= \sum_{k=1}^c \frac{\partial J}{\partial z_k} \cdot \frac{\partial z_k}{\partial \text{net}_k} \cdot \frac{\partial \text{net}_k}{\partial y_j} \\ &= \sum_{k=1}^c -(\bar{z}_k - z_k) \cdot f'(\text{net}_k) \cdot w_{kj} \\ &= \sum_{k=1}^c \delta_k \cdot w_{kj}, \text{ Put this in } ① \end{aligned}$$

$$\begin{aligned} \Rightarrow \frac{\partial J}{\partial w_{ji}} &= \left(\sum_{k=1}^c \delta_k \cdot w_{kj} \right) \cdot f'(\text{net}_j) \cdot x_i \\ &= \delta_j x_i \end{aligned}$$

$$\cdot \Delta w_{ji} = -n \frac{\partial J}{\partial w_{ji}} = n(\delta_j x_i)$$

$$\Delta w_{kj} = -n \frac{\partial J}{\partial w_{kj}} = \cancel{n} (\bar{z}_k - z_k) \cdot f'(\text{net}_k) \cdot y_j$$

* Note: δ_k & δ_j = sensitivity of backpropagation.



; Stochastic BackProp:

- Initialize weights randomly
 - Select samples randomly and modify weights.
- ⇒ At a certain point, test/validation error rises even though training error keeps decreasing, this is when we stop.

—————x————

* Activation Function:

- Should be non linear (Can be linear for small values)
- Should saturate at some point
- Should be continuous and smooth.
- Should be defined
- (Can be monotonically increasing)

a) Linear:

$$y = x$$

b) Sigmoid (Logistic) \Rightarrow 0 to 1 range

$$y = \frac{1}{1 + e^{-x}} = \frac{e^x}{1 + e^x}, \quad \frac{dy}{dx} = \frac{-(-e^{-x})}{(1 + e^{-x})^2} = y - (1-y)$$

c) Hyperbolic Tangent \Rightarrow -1 to 1 range

$$y = a \left[\frac{e^{bx} - e^{-bx}}{e^{bx} + e^{-bx}} \right] \quad (a=1.71, b=2/3) \quad \text{Generally}$$

—————x————

* Practical Aspects:

$$X = [x_1, \dots, x_m]^T_{m \times d}, \quad \tilde{X} = X - \text{mean}(X)$$

- $X_{\text{Norm}} = \tilde{X}/\sigma \rightarrow$ standard deviation of each dimension.
- Use this normalized data only!

- For targets prefer values in range $[-1, 1]$

\Rightarrow Add some random noise to inputs (Training with noise)
(To increase training data size)

- To generate more training data, add translation, rotation etc $\xrightarrow{\text{to}}$ training data.

- * • In general, prefer to have total num of weights to be around $n/10$ (n =Training samples)

\Rightarrow ~~Initialise~~ Initialize weights to random -ve & +ve values, not zeros.

— X —

8 //

- Backprop as feature mapping: Intermediate layer outputs can be used as features.

— X —

1) Learning rates:

- For quadratic error criterion J .

$$\eta_{\text{opt}} = \left(\frac{\delta^2 J}{\delta w^2} \right)^{-1} \quad \text{Newton's Methods can be used as well.}$$

* Momentum:

$$2). w(m+1) = w(m) + (1-\alpha) \Delta w(m) + \alpha \Delta w(m-1)$$

- Use inertia during weight update.

* Weight Decay: (Normalization)

$$3). w_{\text{new}} = w_{\text{old}} (1 - \varepsilon)$$

\Rightarrow Weights which are not updated frequently will go to 0.

* Hints: (Additional nodes in output layer)

- 4). While training, have extra output nodes which help train the network better in backprop.

- These will be removed during testing since we don't require them.

5) Online, Batch or Stochastic training

- Stochastic is preferred over batch
- Online when storing samples is not possible while training.

6) Criterion Functions: (Note: Softmax layer & Logsoftmax loss function at back of book)

- Minkowski norm \rightarrow generalization of Squared sum error
- A few entropy based criterions.

* Convolutional Neural Networks:

- Use multiple filters convolve over entire image.
- Learn filter weights.
(Translational equivariance)

\Rightarrow Spatial Pooling:

- Reduce image size as we move through the network.

* Recurrent Neural Networks: (& LSTM)

- Outputs are fed back into the network as inputs.
- Good for linguistics etc.



Probability & Bayes

$$\bullet P(A|B) = P(A, B)/P(B)$$

→ Prior knowledge based classification:

↳ Consider 2 classes w_1 & w_2

- State of nature = Variable 'w', it can either be w_1 or w_2 . $P(w_1) = x_1, P(w_2) = x_2$

↳ $P(w_1) + P(w_2) = 1$ } Based on prior knowledge

↳ If $P(w_1) > P(w_2)$, then $w \rightarrow w_1$, else $w \rightarrow w_2$.

$$\Rightarrow P(w_j | x) = P(w_j | x) \cdot p(x) = p(x|w_j) \cdot P(w_j)$$

(likelihood) (prior)

* $\Rightarrow P(w_j | x) \propto \frac{p(x|w_j) \cdot P(w_j)}{p(x)}$
 (posterior) (evidence)

$(p(x) = \sum_j p(x|w_j) \cdot P(w_j)$, Just a normalizing term)

$$\Rightarrow P(w_j | x) \propto p(x|w_j) \cdot P(w_j)$$

↳ Now instead of just $P(w_1), P(w_2)$ we use $P(w_1|x)$ & $P(w_2|x)$ (where our sample has x as the feature)

↑
Bayes decision rule

$$\bullet \text{Sensitivity, Recall, Hitrate} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

- We want sensitivity to be as high as possible.

- ex: Cancer test, we want to reduce false negatives

$$\bullet \text{Precision} = \frac{TP}{TP+FP} \quad \bullet \text{Fallout} = \frac{FP}{FP+TN}$$

$$\bullet \text{Specificity} = \frac{TN}{TN+FP} \quad \bullet \text{False Discovery rate} = \frac{FP}{FP+TP}$$

- * Loss $\lambda_{ij} = \lambda(\alpha_i | w_j)$
 - $\Rightarrow \alpha_i = \text{Action taken (Assigned Label)}$
 - $w_j = \text{State of nature (Actual Label)}$
- * Bayes risk = $R(\alpha_i | x) = \sum_{j=1}^c \lambda(\alpha_i | w_j) p(w_j)$
 - We choose α_i if $R(\alpha_i | x) < R(\alpha_j | x)$
- \Rightarrow We often use $\frac{p(x | w_1)}{p(x | w_2)}$ as a decision boundary to select a class
 \nearrow
 We can get these from R , using comparison by bayes rule.

10

— END of Mid —

Expectation:

- * $E[f(x)] = \sum_{x \in D} f(x) \cdot P(x)$

$$E[f(x)] = \int_{-\infty}^{\infty} f(x) \cdot P(x) \cdot dx$$

* Covariance Matrix:

- * $\text{cov}(x, y) = E[(x - E(x))(y - E(y))]$

$$\therefore = E[xy] - E[x]E[y]$$

(n datapoints,
m features)
(Covariance
between features)

$$\Sigma = \begin{bmatrix} \text{cov}(x_1, x_1) & \text{cov}(x_1, x_2) & \dots & \dots & \dots \\ \vdots & & & & \\ & & & \dots & \text{cov}(x_m, x_m) \end{bmatrix}_{m \times m}$$

Covariance
matrix.

$$\rightarrow \text{cov}(x, x) = \text{variance}(x)$$

(Using this we can establish relation b/w different features using multiple samples)

(Note: Eigenvalues
eigen vectors)

* Entropy:

$$H(p(x)) = - \sum_{k=1}^K p(x=k) \log(p(x=k))$$

- Entropy is maximized if $p(x=k) = \frac{1}{K}$
(Uniform distribution)
- Entropy is minimized for dirac-delta function
 $\hookrightarrow p(x=k) = 1$, for one k & 0 elsewhere

* KL Divergence:

- Note: If $y = Ax + b$, $\frac{\partial y}{\partial x} = \text{variance}(y) = \text{var}(Ax+b)$
(Independent of b)

* Normal Densities:

1) Univariate Normal Density:

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

(Gaussian)

- 95% of area is under $\mu - 2\sigma$ to $\mu + 2\sigma$
- Data that is randomly sampled, we can assume it comes from an underlying gaussian.

2) Multivariate Normal Density:

3) Whitening Transform:

- Transform arbitrary set of variables X with a covariance matrix to a new set of variables Y with covariance matrix of $\Sigma = \text{Identity matrix}$.
 \Rightarrow Input vector changes into white noise vector.

$\circ X, \Sigma \rightarrow Y, I$ $\phi^{-1} = \phi$
 $\xrightarrow{\text{single value decomposition theorem}}$

$$\Sigma = \phi \Lambda \phi^{-1} \Rightarrow \text{SVD of } \Sigma \quad , \quad \phi = [\phi_1, \dots, \phi_d]$$

$$A_w = \phi \Lambda^{1/2}, \quad Y = w_0 X \quad \Lambda = \text{diag}(\lambda_1, \dots, \lambda_d)$$

$$A_w^T \Sigma A_w = I, \quad \text{Now } Y = \cancel{A_w^T \Sigma A_w} \cancel{w_0 X}$$

~~level of Σ is λ_i~~

4) Mahalanobis Distance Metric:

$g_1^2 = (x - \mu)^T S^{-1} (x - \mu)$ (Skewed distance)
 $\boxed{g_1^2 = (x - \mu)^T \Sigma^{-1} (x - \mu)}$ } Squared Mahalanobis Distance

$$\Rightarrow g_1^2 = (x - y)^T S^{-1} (x - y) \quad \} \text{Distance b/w 2 random variables}$$

If $S=I$, then this is the euclidean distance
 $(\text{ex: b/w 2 features we have})$

$$\Sigma^{-1} = U \Lambda^{-1} U^T = \sum_{i=1}^d \frac{1}{\lambda_i} u_i u_i^T$$

$$\begin{aligned} \therefore g_1^2 &= (x - \mu)^T \Sigma^{-1} (x - \mu) \\ &= \sum_{i=1}^d \frac{1}{\lambda_i} (x - \mu)^T u_i u_i^T (x - \mu) \\ &= \cancel{\sum_{i=1}^d \frac{1}{\lambda_i}} \cancel{x^T} \cancel{x} \end{aligned}$$

Note: Mid1 \Rightarrow Chapters 1, (2-1-2.3), 4,
(5.1-5.7, 5.8.1, 5.9),
(6.1-6.5, 6.8, 6.10*)

* Multicategory Discriminant Functions:

* Consider $g_i(x)$ for different classes one at a time. Calculate ~~$\hat{w}_i^T x$~~ costs and decide on a class for the input.

$(g_i(x) = w_i^T x) \rightarrow$ Extend this notion.

\Rightarrow But we can use Bayes,

$$g_i(x) = P(w_i|x) = \frac{P(x|w_i) \cdot P(w_i)}{\sum_j P(x|w_j) \cdot P(w_j)}$$

$$\text{or } g_i(x) = P(x|w_i) \cdot P(w_i)$$

$$(g_i(x) = -R(x_i|x)) \Rightarrow \text{we can make this assumption}$$

- Apply a log to $g_i(x)$, and consider this the new $g_i(x)$

$$g(x) \equiv \ln(g_i(x)) = \ln P(x|w_i) + \ln P(w_i)$$

- \rightarrow We generally choose class i if $g_i(x) - g_j(x) > 0$

\therefore This means, $\ln(g_i(x))$

$$\ln(g(x)) = \ln(g_i(x) - g_j(x)) = \ln\left(\frac{P(x|w_i)}{P(x|w_j)}\right) + \ln\left(\frac{P(w_i)}{P(w_j)}\right)$$

(Suppose its 2 class
& we choose class i)

\Rightarrow Suppose we assume $P(x|w_i)$ comes from a normal multivariate density. \rightarrow likelihood density.

$$\begin{aligned} \ln(g_i(x)) &= -\frac{1}{2}(x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) - \frac{d}{2} \ln(2\pi) \\ &\quad - \frac{1}{2} \ln(\Sigma_i) + \ln P(w_i) \end{aligned}$$

- We don't care about constants or terms

Case 1: (Hypspherical Clusters)

- * $\Sigma_i = \sigma^2 I$ } For all classes.

$$|\Sigma_i| = \sigma^{2d}$$

$$\Sigma_i^{-1} = (1/\sigma^2) I$$

- * $\mathbb{I}(g_i(x)) = -\frac{\|x - \mu_i\|^2}{2\sigma^2} + \ln P(w_i)$

(Since we can drop the invariant terms)

($\Sigma_i = \text{constant for all } i$)

\Rightarrow It's similar to distance metric but we also consider prior probabilities.

- * $\mathbb{I}(g_i(x)) = -\frac{(x - \mu_i)^T (x - \mu_i)}{2\sigma^2} + \ln P(w_i)$

$$= -\frac{1}{2\sigma^2} (x^T x + \mu_i^T \mu_i - 2\mu_i^T x)$$

(Since we ignore terms not dependent on class) $+ \ln P(w_i)$

$$= -\frac{1}{2\sigma^2} \mu_i^T \mu_i + \frac{\mu_i^T}{\sigma^2} x + \ln P(w_i)$$

Finally, $\mathbb{I}(g_i(x)) \Rightarrow w_i^T x + w_{i0}$ } This is of the form of $g(x)$ we had initially

- * $w_i = \frac{1}{\sigma^2} \mu_i$

$$w_{i0} = -\frac{1}{2\sigma^2} \mu_i^T \mu_i + \ln P(w_i)$$

- * Now getting back to $\mathbb{I}(g(x)) = \mathbb{I}(g_i(x) - g_j(x))$

$$\begin{aligned}
 g_i(x) &= \left(\frac{\mu_i - \mu_j}{\sigma^2} \right)^T x - \frac{1}{2\sigma^2} \mu_i^T \mu_i + \frac{1}{2\sigma^2} \mu_j^T \mu_j \\
 &\quad + \ln(P(\omega_i)) \\
 &\quad - \ln P(\omega_j) \\
 &= \left(\frac{\mu_i - \mu_j}{\sigma^2} \right)^T x - \frac{1}{2\sigma^2} (\mu_i - \mu_j)^T (\mu_i + \mu_j) \\
 &\quad - \ln(P(\omega_i)/P(\omega_j)) \\
 &= \left(\frac{\mu_i - \mu_j}{\sigma^2} \right)^T \left(x - \left\{ \frac{(\mu_i + \mu_j)}{2} - \frac{\ln(P(\omega_i)/P(\omega_j))}{(\mu_i - \mu_j)^T} \right\} \right)
 \end{aligned}$$

$$g_i(g_i(x)) = w^T (x - x_0) \quad \text{Constant}$$

$g(x) = 0$ (If both classes are equally probable)

* Case 2: (Hyperellipsoidal clusters)

$\Sigma_i = \Sigma$ } for all i , but it's not
Similar to above, we get, $\Sigma \neq \sigma^2 I$

$$\Rightarrow g_i(x) = w_i^T x + w_{i0}$$

$$w_i = \Sigma^{-1} \mu_i$$

$$w_{i0} = -\frac{1}{2} \mu_i^T \Sigma^{-1} \mu_i + \ln P(\omega_i)$$

* Case 3: (Hyperquadrilateral Clusters)

Σ_i = Arbitrary. (we can't

$$g_i(x) = x^T w_i x + w_i^T x + w_{i0}$$

$$w_i = -\frac{1}{2} \Sigma_i^{-1}, w_i = \Sigma_i^{-1} \mu_i$$

$$w_{i0} = -\frac{1}{2} \mu_i^T \Sigma_i^{-1} \mu_i - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i)$$

— x —

* Receiver Operating Characteristic Curves (ROC Curves)

- Consider tasks like object detection, we draw up curves of (Object detected & Object not detected) - (Sensitivity vs Specificity plots.)
- \Rightarrow the curves could be used to detect precision of machine etc.

— x —

12

* Parameter Estimation: (Mean & Variance)

- In the DF's we saw earlier, we assumed likelihood $= P(X|w_i) = l(\theta)$
- \Rightarrow We assumed its a normal distribution, but how do we figure out the parameters?

$D = \{x_1, x_2, \dots, x_n\}$ } Input training data where all these belong to w_i

\Rightarrow We need to estimate μ_i & σ_i :

$$P(x_k | w_i) = N(\theta) \Rightarrow l(\theta) = \ln(p(D|\theta))$$

$\theta = \text{Mean}$
 $\text{we are trying to find out}$

we want to maximize $l(\theta) = \ln(p(D|\theta)) = \ln \prod_{j=1}^n p(x_j | \theta)$

1) Maximum Likelihood Estimation: (MLE)

* a) $\hat{\theta} = \arg \max_{\theta} l(\theta)$

$$p(x_k | \theta) = \frac{1}{(2\pi)^{1/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} (x-\theta)^T \Sigma^{-1} (x-\theta)}$$

$$\ln(p(x_k | \theta)) = \frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln |\Sigma| - \frac{1}{2} (x-\theta)^T \Sigma^{-1} (x-\theta)$$

- Since we want to maximize $\ell(\theta)$ in maximum likelihood estimation,

$$\Rightarrow \nabla_{\theta} \ell(\theta) = \nabla \sum_{k=1}^n \ln p(x_k | \theta) = 0 \quad \left\{ \text{set it to 0} \right.$$

(Log is a monotonically increasing function,
so log & maximize is fine.)

$$\therefore \nabla_{\theta} \ln p(x_k | \theta) = \sum^{-1} (x - \theta)$$

$$\nabla_{\theta} \ell(\theta) = \sum_{k=1}^n (\sum^{-1} (x_k - \theta)) = 0$$

* $\therefore \boxed{\frac{1}{n} \sum_{k=1}^n x_k = \theta}$ \rightarrow This is our mean (as we might have expected)

— x —

- b) Unknown M & $\Sigma = \sigma^2$: (Unlike before where we assumed we knew σ)

* $\ln(p(x_k | \theta)) = -\frac{1}{2} \ln(2\pi\theta_2)$

$$+ -\frac{1}{2\theta_2} (x_k - \theta_1)^2$$

$\Rightarrow \theta_1 = \text{Mean}, \theta_2 = \text{Variance}$

$$\nabla_{\theta} \ell(\theta) = \begin{bmatrix} \nabla_{\theta_1} \ell(\theta) \\ \nabla_{\theta_2} \ell(\theta) \end{bmatrix} \quad \left\{ \text{Both would be 0.} \right.$$

Consider ∇_{θ_2} ,

$$\theta_2 = \nabla_{\theta_2} \ell(\theta) = \nabla_{\theta_2} \left(\sum_{k=1}^n \left(-\frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln \theta_2 - \frac{1}{2} \frac{(x_k - \theta_1)^2}{\theta_2} \right) \right)$$

$$\Rightarrow \sum_{k=1}^n \frac{-1}{2\theta_2} + \frac{1}{2} \frac{(x_k - \theta_1)^2}{\theta_2^2} = 0$$

$$\Rightarrow -n\theta_2 + \sum_{k=1}^n (x_k - \theta_1)^2 = 0$$

$$\left\{ \sigma^2 = \hat{\theta}_2 = \frac{1}{n} \sum_{k=1}^n (x_k - \hat{\theta})^2 \right\}$$

As we got above.
 $\hat{\theta}_1 = \frac{1}{n} \sum_{k=1}^n x_k$

13//

2) Bayesian Parameter Estimation: (BPE)

**

$$P(w_i|x, D) = \frac{p(x|w_i, D_i)p(w_i)}{\sum_{j=1}^c p(x|w_j, D_j)p(w_j)}$$

Class conditional densities, for a certain w_i .

$$\Rightarrow p(x|D) = \int p(x|\theta|D) d\theta = \int p(x|\theta)p(\theta|D)$$

$$\text{Now } p(\theta|D) = \frac{p(D|\theta) \cdot p(\theta)}{\int p(D|\theta) \cdot p(\theta) d\theta} \quad \left. \begin{array}{l} \text{From} \\ \text{Bayes} \end{array} \right\}$$

- We are saying that there are 2 gaussians present, one is of the data & the other is of the means itself \rightarrow It's also a variable curve.

$$\Rightarrow \text{If } p(\theta|D) \text{ peaks at some value, } \hat{\theta}, \text{ then } p(x|D) = p(x|\hat{\theta})$$

$$\text{ie. } p(x|\mu) \Rightarrow \mathcal{N}(\mu, \sigma^2) \quad \left. \begin{array}{l} \\ p(\mu) \Rightarrow \mathcal{N}(\mu_0, \sigma_0^2) \end{array} \right\} \text{Normal distributions}$$

a) Now, $p(\mu|D) = \frac{p(D|\mu) \cdot p(\mu)}{\int p(D|\mu) \cdot p(\mu) d\mu} \quad \left. \begin{array}{l} \\ \text{Constant} \end{array} \right\}$

(Consider univariate gaussians)

$$= \propto \prod_{k=1}^n p(x_k|\mu) \cdot p(\mu)$$

$$= \propto \prod_{k=1}^n e^{-\frac{1}{2} \left(\frac{x_k - \mu}{\sigma} \right)^2} e^{-\frac{1}{2} \left(\frac{\mu - \mu_0}{\sigma_0} \right)^2}$$

$$= \cancel{\propto \prod_{k=1}^n e^{-\frac{1}{2} \left(\frac{(x_k - \mu)^2}{\sigma^2} + \frac{(\mu - \mu_0)^2}{\sigma_0^2} \right)}}$$

$$p(\mu | D) = \propto e^{-\frac{1}{2} \left(\frac{(\mu - \mu_0)^2}{\sigma_0^2} + \sum_{i=1}^n \left(\frac{(x_i - \mu)^2}{\sigma^2} \right) \right)}$$

$$= \propto' \exp \left[-\frac{1}{2} \left(\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \right) \mu^2 - 2\mu \left(\frac{\mu_0}{\sigma_0^2} + \frac{1}{\sigma^2} \sum x_i \right) \right] \rightarrow ①$$

(The other constants not dependent on μ are ignored)

- Now comparing this with the form

Assume it's also a gaussian

$$\therefore \Rightarrow \frac{1}{\sigma_n^2} = \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}$$

$$= \frac{n\sigma_0^2 + \sigma^2}{\sigma^2 \sigma_0^2}$$

$$p(\mu | D) \left\{ \begin{array}{l} * \\ * \\ \text{gaussian's parameters} \end{array} \right\} \Rightarrow \boxed{\mu_n = \left(\frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \right) \bar{x}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \mu_0}$$

- As $n \rightarrow \infty$, $\mu_n = \bar{x}_n$, $\sigma_n \rightarrow 0$

- Mean converges to sample mean
- $\mu_n = \mu_0$ if $\sigma_0 = 0$
- If $\sigma_0 \gg \sigma$, then $\mu_n = \text{sample mean}$.

\Rightarrow We ~~can~~ ^{also} derive that, by putting this in ①

$$p(x|D) \sim N(\mu_n, \sigma_n^2) \quad (\text{Slider})$$

If we know this, we know prob of x belonging to class with sample distribution.

Note: Recursive Bayes Learning,

- Consider ~~a~~ samples, learn parameters
- Add another sample & update

Repeat these

* MLE vs BPE : (Slides)

- MLE is easier to calculate.

⇒ BPE → Bayes or Indistinguishability errors



↳ Overlapping class conditional

densities (The underlying distib.
are of this kind)

→ Model Error

↳ Disparity b/w true / assume
distributio

→ Estimation Error

↳ Small data size.

↳
Not
Gaussi

?



* Principal Component Analysis : (PCA)

- Move from n dimensional data to $k \ll n$ dimensional data so that the dimensions are as independent as possible.

(Try to have minimal correlation b/w dimensions)

→ Preserve the most of the variance of the data after dimensionality reduction.

$$X_{\text{data}} = [x_1, x_2, \dots, x_n] \quad x_i \in \mathbb{R}^d$$

1) • $J_0(x_0) = \sum_{i=1}^n \|x_i - x_0\|^2$, we want to minimize this error

Let $\therefore \nabla J_0(x_0) = 2 \sum_{i=1}^n (x_i - x_0) = 0$,

$x_i = x_0$

(0 Dimensional)

$$\therefore nx_0 = \sum_{i=1}^n x_i$$

$$x_0 = \frac{1}{n} \sum_{i=1}^n x_i \quad \left. \right\} \text{This is the mean}$$

* (1 Dimensional variant) \rightarrow Single eigen vector & value

2) $\hat{x}_i = a_i e + m \quad \left. \right\} \text{Projection of data}$

* (Let $\|e\| = e^T e$) $\quad \text{along some axis}$
 $(\text{Unit vector}) \quad \text{shifted towards mean.}$

$$J_1(a_1, a_2, \dots, a_m, e) = \sum_{i=1}^n \| (a_i e + m) - x_i \|^2$$

$$= \sum_{i=1}^n \| a_i e - (x_i - m) \|^2$$

$$= \sum_i a_i^2 e^T e - 2 a_i e^T (x_i - m)$$

$$+ \| x_i - m \|^2$$

$$J_1(a_1, a_2, \dots, a_m, e) = \sum_i a_i^2 - 2 \sum_i a_i e^T (x_i - m)$$

$$\rightarrow \textcircled{1} + \| x_i - m \|^2$$

$$\nabla_{a_i} J_1 = 2 a_i - 2 e^T (x_i - m) = 0 \quad (\text{Minimize})$$

* $a_i = e^T (x_i - m) \quad \left. \right\} \text{Substitute in } \textcircled{1},$

$$\therefore J_1(e) = \sum_i a_i^2 - 2 \sum_i a_i^2 + \sum \| x_i - m \|^2$$

$$= - \sum_i a_i^2 + \sum \| x_i - m \|^2$$

$$\therefore J_1(e) = - \sum_i e^T (x_i - m) (x_i - m)^T e + \text{const}$$

$$J_1(e) = - e^T S_{d \times d} e, \quad \text{Subject to } \underbrace{e^T e}_{\text{Scatter matrix of data}} = 1$$

$$J_1(e) = - e^T S_e + \lambda (e^T e - 1) \quad \left. \right\} \text{Lagrange optimization}$$

$$\nabla_{\mathbf{e}} J_1(\mathbf{e}) = -2S\mathbf{e} + 2\lambda\mathbf{e} = 0$$

$\therefore S\mathbf{e} = \lambda\mathbf{e}$ } Eigenvectors
} eigenvalues of
 $S \in \mathbb{R}^{d \times d}$

3) (k -dimensional variant)

$$\mathbf{x} = \mathbf{m} + \sum_{i=1}^k a_i \mathbf{e}_i$$

Note: We can prove that $\lambda = \text{Variance}$ for data after
projection $\rightarrow \frac{1}{n} \sum_{i=1}^n \left(\sum_{j=1}^k x_{ij} e_j - \bar{x} \right)^2 = \lambda$. (Video proof)

Note: This PCA can be applied to face classification
(eigen faces)

Note: Once you have matrix of eigen vectors

* $E = [e_1 \ e_2 \ \dots \ e_k] = k \times d$ dimensional

x_i = $d \times 1$ dimensional data point

$\hat{x}_i = E \cdot (x_i - m)$ } we get a k dimensional
data point
Subtract mean

Note: How to choose k ? $\Rightarrow \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^d \lambda_i} = p$
OK
(Sorted order of eigen values)

Try keep $p \in [0.8, 0.95]$

\Rightarrow If we drop more dimensions we lose more, but after a certain point the eigen values become quite small compared to others so we can stop at that point.

15

* Fisher's Linear Discriminant Analysis: (LDA)

- PCA tries to find a projection of data so as to maximize overall scatter.
- LDA tries to find a projection so as to minimize intra-class scatter & maximizing inter-class scatter.

$$\Rightarrow y_i = w^T x_i \quad \left\{ \begin{array}{l} \text{Projected data} \\ \text{Projection of } x_i \text{ data on } w. \end{array} \right.$$

$$X = (\{x_1, \dots, x_{n_1}\}, \{x_{n_1+1}, \dots, x_{n_2}\}) \quad \left\{ \begin{array}{l} \text{2 classes} \\ \text{datapoints} \end{array} \right.$$

$$\begin{aligned} w^T m_1 = \bar{m}_1 &= \frac{1}{n_1} \sum_{i=1}^{n_1} y_i \\ w^T m_2 = \bar{m}_2 &= \frac{1}{n_2} \sum_{i=1}^{n_2} y_{i+n_1} \end{aligned} \quad \left\{ \begin{array}{l} \text{Means of 2 classes.} \\ \text{...} \end{array} \right.$$

$$\Rightarrow S_i^2 = \frac{1}{n_i} \sum_{\substack{y_i \in \text{class} \\ i^{\text{th}} \text{ class}}} (y_i - \bar{m}_i)^2$$

$$\bar{S}_i^2 = \frac{1}{n_i} \sum_{\substack{x \in X_i \\ x \in X}} (w^T x - w^T m_i)^2$$

$$\bar{S}_i^2 = \sum_{x \in X_i} w^T (x - m_i)(x - m_i)^T w$$

$$= w^T S_i w$$

Scatter matrix for a single class

\Rightarrow Now our ~~criterion~~ function,

$$J(\omega) = \frac{|\bar{m}_1 - \bar{m}_2|^2}{\bar{S}_1^2 + \bar{S}_2^2} \quad \begin{cases} \text{Maximize inter class scatter} \\ \text{Minimize intra class scatter} \end{cases}$$

We want to maximize this.

$$= \frac{\omega^T (\bar{m}_1 - \bar{m}_2)(\bar{m}_1 - \bar{m}_2)^T \omega}{\omega^T S_1 \omega + \omega^T S_2 \omega}$$

$$J(\omega) = \frac{\omega^T S_B \omega}{\omega^T (S_1 + S_2) \omega} = \frac{\omega^T S_B \omega}{\omega^T S_w \omega}$$

($S_w = S_1 + S_2$)

$$\Rightarrow \alpha(\omega, \lambda) = \omega^T S_B \omega - \lambda \omega^T S_w \omega$$

$$\nabla_{\omega} \alpha(\omega, \lambda) = 0 = 2S_B \omega - 2\lambda S_w \omega$$

Maximize

$$\Rightarrow S_B \omega = \lambda S_w \omega$$

$$\Rightarrow (S_w^{-1} S_B) \underline{\omega} = \lambda \underline{\omega}$$

Note: Rank of a matrix = Number of non zero eigen values.

~~OK~~ A matrix is invertible if no eigen value is 0 .

So our ω 's are basically eigen vectors of $(S_w^{-1} S_B)$.

\Rightarrow Here we consider a single eigen vector so just one ω corresponding to the maximum eigen value. (We could use multiple eigen vectors as well)

* Multiple Discriminant Analysis (MDA)

$$\cdot S_B = \sum_{i=1}^C N_i (\bar{m}_i - \bar{m})(\bar{m}_i - \bar{m})^T$$

within sample variance

* Generalizes
LDA for
multiple
classes.

$$\cdot S_w = \sum_{i=1}^C \sum_{x \in D_i} (x - \bar{m}_i)(x - \bar{m}_i)^T$$



$$S_w = S_1 + S_2 + \dots + S_C \quad \left. \begin{array}{l} \text{rank is at most } N-C, \\ \text{so it's always singular.} \end{array} \right\}$$

$$\cdot \text{We want to maximize } J(\omega) = \frac{|\omega^T S_B \omega|}{|\omega^T S_w \omega|}$$

$$\Rightarrow S_T \text{ (Total Scatter)} = \sum_{j=1}^N (x_j - \bar{m})(x_j - \bar{m})^T$$

$$S_T = \sum_{i=1}^C \sum_{x \in D_i} (x - \bar{m}_i + \bar{m}_i - \bar{m})(x - \bar{m}_i + \bar{m}_i - \bar{m})^T$$

$$\begin{aligned} ? &= \sum_{i=1}^C \sum_{x \in D_i} (x - \bar{m}_i)(x - \bar{m}_i)^T \\ &\quad + \sum_{i=1}^C N_i (\bar{m}_i - \bar{m})(\bar{m}_i - \bar{m})^T \end{aligned}$$

$$\therefore S_T = S_w + S_B$$



? • Now the issue is that S_w is not invertible

? * \Rightarrow So we apply fisherface solution using PCA
* to reduce dimensionality.

$$* \omega_{PCA} = \max_{\omega} |\omega^T S_T \omega|, S_T = \sum_{j=1}^N (x_j - \bar{m})(x_j - \bar{m})^T$$

\Rightarrow We just replace ~~S_B~~ with $\omega_{PCA} S_B \omega_{PCA}^T$
& S_w with ~~$\omega_{PCA} S_w \omega_{PCA}^T$~~ in $J(\omega)$.

* Fisher Faces: (6 Eigen Faces)

- Similar to eigen faces but instead of eigen vectors from PCA, we use LDA analysis.
- ⇒ Much better than PCA since we optimize on putting ~~the~~ examples of the same class as close as possible & move different classes as far as possible.
- Once we get the LDA vectors, just project the data on them and then use that data to classify the faces.



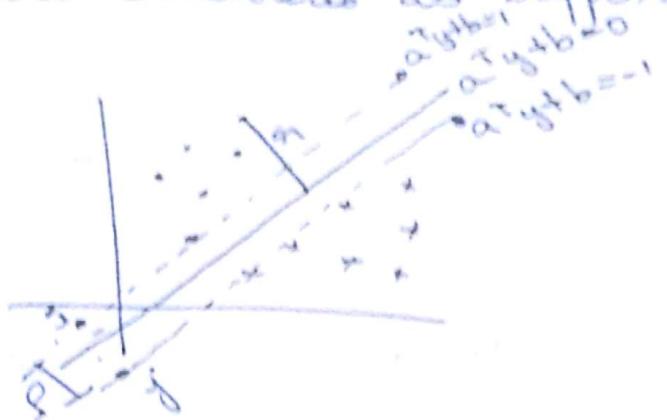
Note: Limitations of Fisher's LDA

- LDA gives at most $c-1$ feature projections.
- LDA assumes unimodal gaussian likelihood.
- LDA fails if class discrimination is captured in the variance & means are similar.
(Since LDA only uses means & doesn't care about the variance)



Support Vector Machine (SVM)

- Data points close to the boundary in all classes are considered as support vectors.



- Now consider a margin around the boundary.

* $g = y_i - y_j \quad \left\{ \begin{array}{l} \text{Assume 2 hypothetical} \\ \text{points on each} \\ \text{of the boundary} \\ \text{with margin} \end{array} \right.$

$$= \frac{a^T y_i + b}{\|a\|} - \frac{a^T y_j + b}{\|a\|}$$

$g = \frac{2}{\|a\|} = \frac{1}{\frac{1}{2} a^T a}$

We want to maximize this

* $a^* = \arg \min_a \frac{1}{2} a^T a \quad \left\{ \begin{array}{l} \text{Primal} \\ \text{Formulation} \end{array} \right.$

(Constraints, $\sum_i (a^T y_i + b) \geq 1 \quad \forall \text{ datapoints}$,
 \downarrow
Class -1 or +1 \rightarrow margin)

- This is a quadratic optimization problem where we optimize a quadratic function subject to linear constraints.

$$* L(a, \alpha, b) = \frac{1}{2} a^T a - \sum \alpha_i (a^T y_i + b z_i)$$

- * This is the dual form where we assign a lagrange multiplier α_i with each constraint.

$$\Rightarrow L(a, \alpha, b) = \frac{1}{2} a^T a - \sum \alpha_i (a^T y_i + b z_i) - \sum \alpha_i b z_i + \sum \alpha_i$$

cong max $\min_{\alpha, b}$ $L(\alpha, a, b)$ such that $a \geq 0$

$$\begin{aligned} \nabla_a L(a, \alpha, b) &= 0 \Rightarrow a = \sum \alpha_i y_i \\ \nabla_b L(a, \alpha, b) &= 0 \Rightarrow \sum \alpha_i z_i = 0 \end{aligned} \quad \left. \begin{array}{l} \text{Substitute these} \\ \text{in (1),} \end{array} \right.$$

$$* \Rightarrow L(\alpha) = \sum \alpha_i - \frac{1}{2} \sum_{i,k} \alpha_i \alpha_k z_i z_k y_i^T y_k$$

Dual (Maximize)
Formulation $\left. \begin{array}{l} \alpha_i \geq 0 \\ \sum \alpha_i z_i = 0 \end{array} \right\}$ constraint.
 $\alpha_i \geq 0$ $\forall i$

$$\frac{\partial L(\alpha)}{\partial \alpha} = 0$$

- This would give us the optimal α_i 's

$$\Rightarrow a = \sum \alpha_i y_i \quad \left. \begin{array}{l} \text{This is very} \\ \text{similar to gradient} \\ \text{descent on page 8} \end{array} \right.$$

$$\Rightarrow b_k = z_k - \sum_{j=1}^n \alpha_j z_j y_j \quad \left. \begin{array}{l} \text{if } \alpha_k > 0 \end{array} \right.$$

- We then take an average of $b_{k=}$ to get the b (margin).

- * Non zero α_k says that y_k is a support vector.

$$\Rightarrow f(y) = \underbrace{\sum_{j=1}^n \alpha_j z_j y_j^\top y + b}_{(a^\top y + b)} \quad \left. \begin{array}{l} \text{The final} \\ \text{classification} \\ \text{function.} \end{array} \right\}$$

Note: This is just like gradient descent but in a statistical way where we get appropriate α_i 's for each y_i . So we set a correspondingly.

$$====x=====$$

* Soft Margin SVM:

- To allow some slack, we try to minimize $\sum \xi_i$ to minimize misclassification additionally

$$z_i(a^\top y_i + b) \geq 1 - \xi_i \quad \left. \begin{array}{l} \text{Set this} \\ \text{additionally} \end{array} \right\}$$

* \Rightarrow So we have, our primal formulation

$$a^* = \arg \min_a \frac{1}{2} a^\top a + C \sum \xi_i$$

$$+ \xi_i \geq 0, \quad z_i(a^\top y_i + b) \geq 1 - \xi_i$$

\hookrightarrow And the dual formulation,

$$\begin{aligned} * L(a, \alpha, b, \beta, \xi) &= \frac{1}{2} a^\top a - \\ &- \sum \alpha_i \{ (a^\top y_i + b) z_i - 1 + \xi_i \} \\ &- \sum \beta_i \xi_i + C \sum \xi_i \end{aligned}$$

$$\begin{aligned} L(a, \alpha, b, \beta, \xi) &= \frac{1}{2} a^\top a - \sum \alpha_i a^\top y_i z_i \\ &- \sum \alpha_i z_i + \sum \alpha_i - \sum \alpha_i \xi_i \\ &- \sum \beta_i \xi_i + C \sum \xi_i - \sum \xi_i (\alpha_i \beta_i) \end{aligned}$$

Now, we make

$$\nabla_{\alpha} L = 0, \nabla_{\beta} L = 0, \nabla_{\gamma} L = 0$$

$$\begin{aligned}\downarrow & \quad \downarrow \\ \alpha_i + \beta_i &= c \quad \alpha = \sum_i \alpha_i y_i \\ 0 \leq \alpha_i &< c\end{aligned}$$

$$\sum_i \alpha_i y_i = 0$$

$$\Rightarrow L = \sum_i \alpha_i - \frac{1}{2} \sum_{i,k} \alpha_i \alpha_k y_i y_k z_i z_k$$

- * We just have an additional constraint here, $\rightarrow 0 \leq \alpha_i < c$.

- We have a similar formulation as before.

===== X =====

17 * Transductive SVM:

- We learn from training & test data simultaneously
 \rightarrow This is a weak semi-supervised setup

\rightarrow Normal SVMs till now were Induction based \rightarrow fully supervised.

- Consider x_1, \dots, x_n $\{$ training data $\}$
 x_{n+1}, \dots, x_m $\{$ Test data $\}$

$$\rightarrow \text{arg min}_{\alpha_{1:n}, \alpha_{n+1:m}} \left(\frac{1}{2} \alpha_{1:n}^T \alpha_{1:n} + C \sum_{i=1}^n \alpha_i \right) + D \sum_{i=n+1}^m \alpha_i$$

Constraints:
 $z_i(\alpha^T x_i + b) \geq 1 - \xi_i, \xi_i \geq 0$ (infeasible)
 $z_i(\alpha^T x_i + b) \leq 1 + \eta_i, \eta_i \geq 0$ (infeasible)

→ We do this iteratively

- ↳ Assume z_{m+1}, \dots, z_m and then minimize others
- ↳ Assume others are constant & minimize z_{m+1}, \dots, z_m

- Repeat those 2 steps.

— — — x — — —

* Multi Category SVM:

- ↳ One vs All $\xrightarrow{\text{other classes}} n$ SVMs
- ↳ One vs One $\xrightarrow{\text{each other class}} n \times (n-1)$ SVMs

$$a_i^* y_k - a_j^* y_k > 0$$

⇒

if $a_i^* y_k - a_j^* y_k > 0$ then y_k is correctly classified by i

if $a_i^* y_k - a_j^* y_k \leq 0$ then y_k is misclassified by i



Non Linear SVM:

- Linear classification in non linear higher dimensional space.

$$\text{arg max}_{\alpha_1, \dots, \alpha_n} \sum_{k=1}^n \alpha_k - \sum_{k \neq j} \alpha_k \alpha_j y_k y_j \phi(x_k)^T \phi(x_j)$$

→ ϕ = Mapping of data point to another space.

* Kernels:

- * Kernel or Gram Matrix

$$K(k, j) = \phi(y_k)^T \cdot \phi(y_j)$$

\downarrow
 $n \times n$
matrix

$$= \langle \phi(y_k), \phi(y_j) \rangle$$

Inner product

$\Rightarrow K \Rightarrow$ positive semi-definite matrix } Eigenvalue ≥ 0

↳ Symmetric

Desired
but not always true

- Any matrix satisfying symmetric & positive semi-definite matrix behaviour is a kernel

i) Linear Kernel

$$K(k, j) = y_k^T y_j$$

ii) Polynomial Kernel

$$K(k, j) = (1 + y_k^T y_j)^p$$

* iii) Gaussian | Radial Basis Function (RBF) Kernel

$$K(k, j) = \exp\left(-\frac{\|y_k - y_j\|^2}{2\sigma^2}\right) = \exp\left(-\gamma \|y_k - y_j\|^2\right)$$

iv) Sigmoid Kernel

$$K(k, j) = \tanh(\beta_0 y_k^T y_j + \beta_1)$$

(This is not positive semi-definite)

- If ϕ is overcomplicated, we can end up overfitting.

→ With kernels we can avoid having to define ϕ explicitly (Useful at times)

* Kernel SVM: (example at the back)

- Only consider support vectors \rightarrow So efficient where there is a lot of data
 - Overfitting can be controlled by soft margin approach
 - Flexibility to choose any kernel.
- ~~→~~ One big issue is outliers, if they get selected as support vectors, we end up detecting completely wrong boundaries.
- Using just the kernel matrix we can find out the SVM parameters.

18/11

* Kernel Trick: (Project lower dimensional data to higher dimensions or infinite dimensional data to finite dimensions)

Note: $K = X^T X \quad \{ \text{Kernel Matrix}$

* $K = U \Lambda U^T \quad \{ \text{Eigen decomp}$

$$\begin{aligned} &= U \Lambda^{1/2} \cdot \Lambda^{1/2} U^T \\ &= U \Lambda^{1/2} \cdot (U^T \Lambda^{1/2})^T \\ &= \phi^T \phi, \quad Y = U^T \Delta^{1/2} \end{aligned}$$

↑
Every kernel should satisfy this

$$\phi: X \rightarrow Y$$

- At times input data could be text etc, (i.e. input dimension is variable & large), we often use various functions to compute the kernel matrix

↳ ex: Edit Distance or Num of common substrings

* For strings.

↳ Kernels for documents

↳ Bag of words, TFIDF

↳ Kernels for graphs

↳ Counting matching random walks

- We can join up multiple kernels one after another and prove they are also kernels.
(Slides for example)

————— * —————

* Kernelized KNN:

- Instead of $\|x_i - x_j\|^2$,

$$\|\phi(x_i) - \phi(x_j)\|^2 = 2K(x_i, x_j) + K(x_i, x_i) + K(x_j, x_j)$$

————— * —————

* Kernelized PCA: (Slides also has the pdf)

- $x \in \mathbb{R}^d$, $\rightarrow S$ a $d \times d$ scatter matrix

$$\Rightarrow Y = \phi(X), Y_i = \phi(x_i) \in \mathbb{R}^{q \times p}$$

$$S' = \sum_{i=1}^n (\phi(x_i) - m)(\phi(x_i) - m)^T \quad \left\{ \text{New scatter matrix} \right.$$

$$\Rightarrow S' = \sum_{i=1}^n \phi(x_i) \phi(x_i)^T \rightarrow \textcircled{1} \quad \begin{array}{l} (\text{Assume } m = 0) \\ (\text{Centered data}) \end{array}$$

- We don't know this scales,
- Replacing $\textcircled{1}$ in above, material since we don't know ' ϕ ' function)
- $S' w_{pi} = \lambda w_{pi}$ } Our new eigen vectors w_i with eigen values λ .
- We want to find these w_i ($w^T w = 1, \forall i$)

$$\frac{1}{\lambda} \sum_i \phi(x_i) \cdot \phi(x_i)^T \omega = \omega \rightarrow \textcircled{2}$$

$$\Rightarrow \sum_i \phi(x_i) \cdot \frac{\phi(x_i)^T \omega}{\lambda} = \omega$$

$$\Rightarrow \sum_i \phi(x_i) \cdot \alpha_i = \omega \quad , \quad \alpha_i = \begin{matrix} \text{Scalar} \\ \text{value} \end{matrix} \rightarrow \textcircled{3}$$

- Substitute \textcircled{3} in \textcircled{2}
- $\left(\sum_i \phi(x_i) \cdot \phi(x_i)^T \right) \sum_k \phi(x_k) \cdot \alpha_k = \lambda \sum_k \phi(x_k) \cdot \alpha_k$

\Rightarrow Multiply both sides by $\phi(x_k)^T + \lambda$

- $\sum_{i,k,l} \phi(x_i)^T \phi(x_i) \cdot \phi(x_i)^T \cdot \phi(x_k) \alpha_k$
 $= \lambda \sum_{k,l} \phi(x_k)^T \phi(x_l) \alpha_k$

$$\Rightarrow K \cdot K \cdot \alpha = \lambda \cdot K \cdot \alpha$$

* Round matrix (we know this) $K \cdot \alpha = \lambda \cdot \alpha \rightarrow \alpha^T \alpha = 1$ we learn \alpha's using kernel matrix

From \textcircled{3},

- Now, $\omega_{\text{pri}} = Y_{\text{pri}} \cdot \alpha_{\text{pri}}$ } we then get ω using our new \alpha's.

All inputs $\phi(X)$ $d \times n$

Now,

$$\omega^T \omega = \alpha^T Y^T Y \alpha$$

$$= \alpha^T K \alpha \rightarrow \lambda \alpha^T K = \lambda$$

We want this

$$\rightarrow \frac{\omega^T \omega}{\lambda} = 1, \quad \hat{\omega} = \frac{Y \alpha}{\sqrt{\lambda}} \quad \left. \begin{matrix} \text{These are} \\ \text{the } \omega's \\ \text{we need.} \end{matrix} \right\}$$

$$\hat{\omega}^T \hat{\omega} = 1$$

Note: The Kernelized methods are very useful since a lot of times input data might not be vectorizable, strings, graphs etc. We define a kernel based on some comparison and then using that we can get vectorizable data.

— X —

Note: Kernel LDA, Multiple Kernel Learning,
* Non-linear dimensionality reduction.

— X —

Note: Mid 2 — 2.1-2.3, 2.5, 2.6, 28.3

(Duda book) 3.1, 3.2, 3.3, 3.4, 3.5.1, 3.7, 3.8
5.11, 5.12

— X —

Data Clustering:

(Unsupervised, no labels given)

Similarity Measures (between datapoints)

- Cosine distance
- Jaccard distance (Set)
- Minkowski distance
(City-block & Euclidean)
- Mahalanobis distance

- Normalization is not always helpful. (Context dependent)

* Criterions: (Minimization approaches)

1) Sum of squared error. (Not always the best criterion)

$$\hookrightarrow J_e = \sum_{j=1}^c \sum_{x \in D_i} \|x - m_i\|^2 , m_i = \text{Mean of } i\text{-th cluster}$$

↑
For each cluster

- Minimizing this via optimization is good but not always the best. (example in slides)

2) Related minimum variance:

$$\star J_e = \frac{1}{2} \sum_{i=1}^c n_i \bar{s}_i , \bar{s}_i = \frac{1}{n_i^2} \sum_{x \in D_i} \sum_{x' \in D_i} \|x - x'\|^2$$

↑
Minimize.

$$(s_i = \frac{1}{n_i^2} \sum_{x \in D_i} \sum_{x' \in D_i} s(x, x') , s_i = \min_{x \in D_i, x' \in D_i} s(x, x'))$$

— x —

* Hierarchical Clustering: (Agglomerative)

- * • Agglomerative clustering (non-parametric model)
(non generative)
1) Bottom Up

⇒ Initially assume all datapoints belong to different clusters.

↳ Keep merging the nearest 2 clusters
(Keep doing this till you get required k clusters)

Cluster Distance measures

$$d_{\min}(D_i, D_j) = \min_{\substack{x \in D_i \\ x' \in D_j}} \|x - x'\|$$

$$d_{avg}(D_i, D_j) = \frac{1}{n_i n_j} \sum_{x \in D_i} \sum_{x' \in D_j} \|x - x'\|$$

$$d_{mean}(D_i, D_j) = \|m_i - m_j\|$$

⇒ Using $d_{\min}(D_i, D_j)$ for merging of clusters by linking the 2 points yields an MST (Minimum Spanning Tree) solution.

- Using d_{max} → We end up getting fully connected subgraphs (large clusters)

⇒ For d_{avg} , we can use any other similarity measure we like.

- d_{mean} = Computationally simple.



2) K-Means Clustering

(Iterative method) (example in slides).

- Initialize means randomly

- Then iterate between

E-step \rightarrow Assign each datapoint to nearest mean.
 M-step \rightarrow Update cluster means.

$$\Rightarrow J_e = \sum_{i=1}^n \sum_{j=1}^k a_{ij} \cdot \|x_i - \mu_j\|^2, a_{ij} \in \{0,1\}$$

(For hard clustering)

$$\left(\sum_{j=1}^k a_{ij} = 1 + i \right)$$

(Each point is assigned to only one cluster)

- We assume that each of the k clusters have same variance so that we can move points into clusters with the nearest mean. So our criterion/error function is correct.

$$\Rightarrow \nabla_{\mu_j} J_e = 0, \text{ To minimize } J_e.$$

$$\star \sum_{i=1}^n a_{ij} \cdot (x_i - \mu_j) = 0 + j \in \{1, k\}$$

~~$\Rightarrow \mu_j = \frac{1}{n} \sum_{i=1}^n x_i$~~ subject to the fact that a_{ij} is known

$$\mu_j = \frac{\sum_{i=1}^n x_i \cdot a_{ij}}{\sum_{i=1}^n a_{ij}}$$

- Similarly if you fix μ_j 's, then we can optimize on a_{ij} .

\Rightarrow This is convergent to local minima \rightarrow so sensitive to initialization. Sensitive to outliers.

- We also always assume our clusters to be convex, so fails for certain distributions.
- Complexity is also quite high $O(N \cdot K \cdot d \cdot L)$.

* a) Probabilistic Kmeans:

- * • Instead of hard classification, we assume a probabilistic assignment of datapoints to clusters.

$0 \leq a_{ik} \leq 1$. (Instead of only 2 possible values)
 (This is tough so we use another formulation)

⇒ We represent probability mixture of data as a gaussian mixture model.

$$p(x) = \sum_{k=1}^K \pi_k \cdot N_k(x; \mu_k, \Sigma_k), \quad \sum_{k=1}^K \pi_k = 1 \quad (\text{Multivariate gaussians}) \quad 0 \leq \pi_k$$

Given all μ_k & Σ_k .

$$p(x) = \sum_{k=1}^K p(k) \cdot p(x|k)$$

Note: K can be determined using Bayesian mixture model (The number of clusters)

⇒ Here it's similar to MLE except that we want $p(x)$ to be maximized.
 $(\log(p(x))) \rightarrow \text{Maximize}$

$$\begin{aligned} * \cdot l(\theta) &= \ln(p(x|\theta)) \\ &= \ln\left(\sum_H p(x, H|\theta)\right) \\ &= \ln\left(\sum_H \frac{q(H)}{q(H)} \cdot p(x, H|\theta)\right) \end{aligned}$$

• We intro
H as a
missing
para

⇒ log of summation is hard to compute

⇒ Jensen's Inequality:

$$\log E(x) \geq E[\log x]$$

Using Jensen's equality,

$$l(\theta) \geq \sum_{\text{H}} q(H) \ln \left(\frac{p(x, H | \theta)}{q(H)} \right)$$

$$l(\theta) \geq \sum_{\text{H}} p(H | x, \theta) \cdot \ln \left(\frac{p(x, H | \theta)}{p(H | x, \theta)} \right)$$

$$l(\theta) \geq \sum_{\text{H}} p(H | x, \theta) \cdot \ln(p(x | \theta))$$

↑
Basically Maximize RHS
so that the lower bound of $l(\theta)$ increases.

Let, $q(H) = p(H | x, \theta)$
 $q(H) = \frac{p(H, x | \theta)}{p(x | \theta)}$

GrMM

- Given datapoints, we try to optimize and achieve M_j 's, Σ_j 's and π_j 's.
→ This is done using what we derived above, optimizing the RHS sum of logs would optimize our $l(\theta)$ which is what we want.

— x —

b) Fuzzy K-Means: (Slides)

- * • Points belong to a cluster in a fuzzy manner. Points near cluster boundary belong lesser to the cluster etc.

C) Kernel K-Means:



- Useful when data has ∞ dimensions etc where we cannot compute the means.

$$J = \sum_{i=1}^N \sum_{j=1}^K a_{ij} \|\phi(x_i) - \tilde{m}_j\|^2$$

$$\tilde{m}_j = \frac{\sum_{i=1}^N a_{ij} \phi(x_i)}{\sum_{i=1}^N a_{ij}}$$

\Rightarrow We move onto kernelized implementation so that we know $\phi(x_i) \cdot \phi(x_j)$ but not ϕ individually. (Derivation back of notes)

Note: Read up on Mixture Models \rightarrow Imp, for example

*3) Spectral Clustering:

- Apply graph based analysis to move classes into blobs so we can use k-means directly. (Projection onto another Euclidean space)

\Rightarrow Initially form a graph from all data points and perform graph partitioning.

- Graph construction \rightarrow k-nearest neighbor from datapoints.
link them
(or ϵ neighborhood
(sphere of all neighbors within))

(Note: Semisupervised methods

also specify if 2 points must belong to a class or must not belong)

a) MinCut Graph Partitioning:

⇒ 2 Way partitioning (2 clusters)

*

- We want to divide it into clusters such that we minimize these

* ~~RatioCut~~ (C_1, \dots, C_k) = $\frac{1}{2} \sum_{i=1}^k \frac{\text{cut}(C_i, \bar{C}_i)}{|C_i|}$

* $\rightarrow \text{NormCut}(C_1, \dots, C_k) = \frac{1}{2} \sum_{i=1}^k \frac{\text{cut}(C_i, \bar{C}_i)}{\text{val}(C_i)}$

(\bar{C}_i = All other points apart from cluster C_i)

- This is NP, so we can't solve it. We move on to a relaxed version

* b) Graph Projection from Partitioning:

- * • We define a function f , which gives a value for each node in the graph.

$$f(v_i) = \begin{cases} \sqrt{|C_i|/\bar{C}} & \text{if } v_i \in C \\ -\sqrt{|C_i|/\bar{C}} & \text{if } v_i \in \bar{C} \end{cases} \quad \begin{array}{l} \text{we will prove} \\ \text{why this is} \\ \text{the right thing} \\ \text{next.} \end{array}$$

⇒ So the ratio cut criterion here is $(k=2)$

~~RatioCut~~ \Rightarrow

$$\Rightarrow |V| \cdot \text{RatioCut}(C, \bar{C}) \Rightarrow (|C| + |\bar{C}|) \left\{ \frac{1}{2} \frac{\text{MinCut}(C, \bar{C})}{|C|} + \frac{1}{2} \frac{\text{MinCut}(C, \bar{C})}{|\bar{C}|} \right\}$$

(we want to minimize this)

$$\Rightarrow \frac{1}{2} \text{Min cut}(c, \bar{c}) \left\{ \frac{|c| + |\bar{c}|}{|c|} + \frac{|c| + |\bar{c}|}{|\bar{c}|} \right\}$$

$$\Rightarrow \frac{1}{2} \text{Min cut}(c, \bar{c}) \left\{ \frac{|\bar{c}|}{|c|} + \frac{|c|}{|\bar{c}|} + 2 \right\}$$

$$\Rightarrow \frac{1}{2} \sum_{\substack{i \in c \\ j \in \bar{c}}} w_{ij} \left(\sqrt{\frac{|\bar{c}|}{|c|}} + \sqrt{\frac{|c|}{|\bar{c}|}} \right)^2 + \frac{1}{2} \sum_{\substack{i \in c \\ j \in \bar{c}}} w_{ij} \left(-\sqrt{\frac{|c|}{|\bar{c}|}} - \sqrt{\frac{|\bar{c}|}{|c|}} \right)$$

Minimise $\frac{1}{2} \sum_{i,j=1}^n w_{ij} (f(v_i) - f(v_j))^2$ } we want to minimise this
w.r.t f . $\rightarrow (1)$ (Quadratic form)

- Retrieving the exact f is also an NP problem so we relax it and try to solve
- ∴ We have showed that $f(v_i)$ is what we assumed it should be initially.

$$\underbrace{\sum_{i=1}^n f(v_i)}_{f \perp 1 \text{ ie orthogonal to } [1, 1, \dots, 1]} = 0, \|f\| = \sqrt{\sum_{i=1}^n f(v_i)^2} = \|f\| = \sqrt{|c| + |\bar{c}|} = \sqrt{n}$$

\Rightarrow Graph Incidence Matrix:

$$*\ *\ \nabla_{ij} = \begin{cases} \nabla_{ev} = -1 & \text{if } v \text{ is initial vertex of edge } ij \\ \nabla_{ev} = 1 & \text{if } v \text{ is terminal vertex of edge } ij \\ \nabla_{ev} = 0 & \text{otherwise} \end{cases}$$

(Each row corresponds to an edge)
(Sum of each row = 0)

\Rightarrow This matrix is like the derivative.

$$(\nabla f)(e_{ij}) = f(v_j) - f(v_i)$$

$$L = \nabla^T \nabla \quad \} \text{Laplace (2nd derivative)}$$

$$\rightarrow (L\mathbf{f})(v_i) = \sum_{v_j \sim v_i} (d(v_i) - f(v_j))$$

$L = D - A$, D = Degree matrix, $d_i = d(v_i)$

- Diagonal entries of L are degree of vertices and the other entries represent the adjacency matrix.

$$(L\mathbf{f})(v_i) = \sum_{v_j \sim v_i} w_{ij}(f(v_i) - f(v_j))$$

Quadratic form, (Minimize this) (from ①)

$$\mathbf{f}^T L \mathbf{f} = \frac{1}{2} \sum_{e_{ij}} w_{ij}(f(v_i) - f(v_j))^2$$

(L = true semi-definite & symmetric)

Finally, we want to optimize this.

$$\Rightarrow \arg \min_{\mathbf{f}} \frac{\mathbf{f}^T L \mathbf{f}}{\mathbf{f}^T \mathbf{f}}$$

- We now know that this is nothing but requiring us to find eigen vectors & values of L .

($\lambda_{\min} = 0$, (the smallest λ) for a connected graph)

\Rightarrow If we want to project to k dimensions
 * we use multiple eigen vectors of L .
 (Non-zero corresponding eigenvalues)

- (Note: If we have multiple connected components in graph, $L = \begin{bmatrix} L_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & L_n \end{bmatrix}$ where each L_i is a Laplacian matrix for that connected component i .)

\Rightarrow Once we project these points, we can apply simple clusterings like k-means etc to get good clustering.

k

Decision Trees

(Classification)

- Nominal data (categorical) is used. We don't use vectorial data.
- Very useful in businesses etc since they can provide stepwise explanation for the outcome we chose.

⇒ { Decision trees → Target is continuous (Regression tree)
* ↳ Target is discrete (Classification tree)
(examples in slides)

- Leaf nodes of tree corresponds to target value.
* (Decision nodes specify test on a single attribute)

⇒ We stop splitting nodes when

- *
 - ↳ i) All elements of training at this point at the node belong to the same class.
 - ↳ ii) We have used up all attributes (We can also use pairs of attributes etc if we want!) (we also need to take care that we don't overfit)

- Every decision tree can be represented as a binary tree. We test ($\text{attr} = \text{val} ?$), and branch as yes/no. These are CART trees (Classification & Regression Trees).

Note: If we choose single attributes in each test,
* then we end up with monothetic trees & boundaries are just multiple hyperplanes, not generalized enough.

* Classification & Regression Trees: (CART)

- * Purity of node = Similarity of target values of training data at that node

⇒ To split a node, we choose an attribute which gives us higher purity in descendants.

i) Entropy Impurity = $-\sum_j P(w_j) \cdot \log_2 P(w_j)$

* (Entropy=0 for a pure node)

ii) Gini Impurity = $\sum_{i \neq j} P(w_i) \cdot P(w_j) = \frac{1}{2} - \sum_j \frac{P^2(w_j)}{2}$

iii) Misclassification Impurity = $1 - \max_j P(w_j)$

- We can also say we choose an attribute which maximizes the decrease in impurity ($\Delta i(N) = i(N) - P_L i(N_L) - (1-P_L) i(N_R)$)

$$\Delta i(N) = i(N) - P_L i(N_L) - (1-P_L) i(N_R)$$

- This is a greedy approach & finds local optima.
- Having branching factor > 2 decreases the impurity more. So we tend to use gain ratios as well (slides)

⇒ We stop splitting based on thresholds, the amount of information gain by best attributes split or the amount of purity or height of tree etc.

* ex: Global Criterion = $\alpha \cdot \text{size}_{\text{tree}} + \sum_{\text{leaf nodes}} i(N)$

* Pruning:

- When we stop splitting too early, it's called the

* Horizon Effect

- In order to account for this, using the global criterion won't work, so what we do is build a complete tree (overfit) and then start merging nodes based on criterions.
(~~the~~ p.m.)
- This is computationally expensive.
- The merging of nodes provides good interpretability.

⇒ We then need to assign labels for impure leaf nodes. Majority voting does this.

- We perform PCA first & then use them in decision trees so that the trees are simple and don't overfit (example in slides)
(we don't want correlation in data)
- At nodes we can use multi attribute functions

⇒ Handling class priors etc can be done by modifying impurity formulae,

$$\text{ex: } i(N) = \sum_{ij} \lambda_{ij} P(w_i) P(w_j)$$

* Missing Attributes: (ex: Slides)

- At a node if primary attribute is missing we define surrogate splits by considering other attributes such that the correlation with primary attribute split is the highest.



Limitations:

- Finding best decision tree is an NP-complete problem.
- Greedy just gives us a local optima.
-

→ handles missing values

→ handles cardinality



→ handles discrete attributes

Note:

- * • ID3 tree (a kind of decision tree)
 - Branches from a node = No. of values per attribute
 - Real valued attributes are binned.
 - Stop when leaf nodes are completely pure.
- C4.5
 - ID3 for categorical attributes
 - CART for real valued attributes.
- ? • Bagging, Random forest, Rotation forest



(Decision boundary)

Attribute \rightarrow A_1, A_2, \dots, A_n

Decision boundary \rightarrow D_1, D_2, \dots, D_m

• Binary

Decision boundary \rightarrow D_1, D_2, \dots, D_m

* Bayesian Nets:

- Graphical Model = Graph Theory + Probability Theory.
 - * \rightarrow Directed \rightarrow Bayesian Nets
 - \nwarrow Undirected \rightarrow MRFs
- $(A) \rightarrow (B)$ indicates that B is caused by A.
 $P(B, A) = P(B|A) \cdot P(A)$
- Top-Down {Reasoning
Bottom-UP {Diagnostic

\Rightarrow Conditional Independence:

- * i) $(A) \rightarrow (B) \rightarrow (C)$ (Head to tail)
 - C is conditionally independent on A given B.
- ii) $\begin{array}{c} \cancel{(A)} \\ \cancel{\longrightarrow} \end{array} \longrightarrow (C)$ (Tail to tail)
 $(B) \longrightarrow (A)$
 - C is conditionally independent on A given B.



- If C is not observed, A, B are conditionally independent. But if C is observed then A depends on B, C and B depends on A, C (Explaining Away)

⇒ Arc Reversal: (Examples in slides)

- Based on Bayes

$$\begin{aligned} p(x_1, x_2) &= p(x_1 | x_2) \cdot p(x_2) \\ &= p(x_2 | x_1) \cdot p(x_1) \end{aligned}$$

(we flipped an edge)

⇒ we could also add a new edge and flip another.

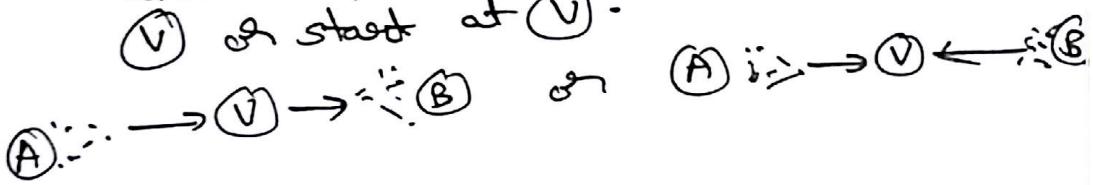
— x —

* D-Separation: (Determine if two variables are conditionally independent)

*

- Two variables (A) & (B) are conditionally independent if,

↳ There is a node (V) which is observed and all paths b/w A & B go through (V) or start at (V) .



— x —

* Note: Given a network and a few nodes (Message where data is observed, we calculate passing) the probabilities of all other nodes by moving to nearby nodes of these nodes. (BFS like)

(we can either move downwards reasoning or move upwards diagnostic to get probabilities of all unknown nodes)

— x —