

Linear Regression:

- Normal eqns, L_2/L_1 regularization
- Auto Regression models
- Convex linear/logistic objective

Classification:

K-NN

- logistic regression; convex objective

Optimisation:

- Gradient descent, momentum, nesterov
- Newton
- Conjugate gradients

Neural Networks:

- Autoencoder, Autoencoder, PCA
- LSTM
- Initialisation

Auto diff:

- Forward; complex & dual
- Backward
- Backprop & parameter tying

Dimension Reduction:

- PCA; eigendecomposition & SVD, ICA
- Fisher's LDA

Clustering:

- K-means
- Mixture Models; GMM
- Spectral clustering

Fast Nearest Neighbours:

- distance metric
- Orchard, AESA
- boids
- k-d tree

Linear Regression

$$E(\underline{w}) = \sum_i (y_i - \underline{w}^T \underline{x}_i)^2 \leftarrow \text{convex as } \underline{H} \text{ positive definite} \\ (\underline{z}^T \underline{H} \underline{z} > 0)$$

$$O(D^3) \text{ time} \rightarrow \underline{\underline{X}}^T \underline{\underline{X}} \underline{w} = \underline{\underline{X}}^T \underline{y}$$

$$\hat{\underline{w}}_{L_2} = (\underline{\underline{X}}^T \underline{\underline{X}} + \lambda \underline{\underline{I}}) \underline{\underline{X}}^T \underline{y} \quad + \text{closed form solution}$$

- not space; try L_1

Auto-regressive models:

• time series V_1, V_2, \dots, V_T

• predict V_T using $V_{T-L+1:T} = \underline{\underline{\Phi}}_T$

• can use least squares / L_1/L_2 regularisation

$$E(\underline{\alpha}) = \sum_t (V_t - \underline{\underline{\Phi}}_t^T \underline{\alpha})^2$$

Classification

• R-NN: Vote of k closest using some distance metric

• Logistic regression:

$$P(C=1 | \underline{x}) = \sigma(\underline{w}^T \underline{x}), \quad \sigma(x) = \frac{e^x}{1+e^x}, \quad \sigma(-x) = 1 - \sigma(x), \\ \sigma'(x) = \sigma(x)(1 - \sigma(x))$$

$$L(\underline{w}) = \sum_i \log P(C_i | \underline{x}_i)$$

$$= \sum_i \left(C_i \log \sigma(\underline{x}_i^T \underline{w}) + (1 - C_i) \log (1 - \sigma(\underline{x}_i^T \underline{w})) \right)$$

• optimise \underline{w} w/ gradient descent methods

Optimisation

$O(DN)$ time
 $O(DN)$ space

• Gradient Descent: $\underline{w}_{k+1} = \underline{w}_k - \epsilon \nabla E(\underline{w})|_{\underline{w}_k}$

• error falls as $1/T$

• Momentum: $\underline{w}_{k+1} = \underline{w}_k + \tilde{\underline{g}}_{k+1}, \quad \tilde{\underline{g}}_{k+1} = \mu_k \tilde{\underline{g}}_k - \epsilon \nabla E(\underline{w})|_{\underline{w}_k}$

• moving average gradient

• Nesterov:

• momentum, but 'smoother' gradient MA update

$$\tilde{\underline{g}}_{k+1} = \mu_k \tilde{\underline{g}}_k - \epsilon \nabla E(\underline{w})|_{\underline{w}_k + \mu_k \tilde{\underline{g}}_k}$$

Simple, but
second order
system dependent

$O(D^2)$
per iteration

• Newton: $\underline{w}_{k+1} = \underline{w}_k - \epsilon \underline{H}_{E(\underline{w})}^{-1} \nabla E(\underline{w}) \Big|_{\underline{w}_k}$

- invariant under coord transform, but \underline{H} expensive to comput & invert ($O(D^3)$)
- 1 step convergence for quadratic $E(\underline{w})$

• Conjugate gradients:

$$p_i \triangleq p_i^\top = \delta_{ij}, p_i \triangleq p_i^\top \Rightarrow p_1, p_2, \dots \text{conjugate vectors for}$$

can do analytically
for quadratic
objective

$$R = 1$$

$$\underline{w}_R = \dots, \underline{p}_R = -\nabla E(\underline{w}) \Big|_{\underline{w}_R}$$

$$\text{while } \nabla E(\underline{w}) \Big|_{\underline{w}_k} \neq 0:$$

line search along $\underline{p}_k \longrightarrow \alpha_k = \underset{\alpha}{\operatorname{argmin}} E(\underline{w} + \alpha \underline{p}_k)$

$$\underline{w}_{k+1} = \underline{w}_k + \alpha_k \underline{p}_k \quad \text{update } \underline{w}$$

→ find new conjugate gradient $\longrightarrow \underline{p}_{k+1} = -\nabla E(\underline{w}) \Big|_{\underline{w}_{k+1}} + \beta_k \underline{p}_k$

• Max D iterations

Neural Networks

$$\underline{x} \rightarrow \underline{h}_1 \rightarrow \underline{h}_2 \rightarrow \dots \rightarrow \underline{y}$$

$$\cdot y = \sum \sigma_j(x \mid \underline{w}) = \sigma_L(\underline{w}_L, \underline{h}_{L-1}), \quad \underline{h}_1 = \sigma_1(\underline{w}_1, \underline{h}_{0-1}), \quad \underline{h}_i = \sigma_i(\underline{w}_i, \underline{x})$$

• optimise by gradient descent

• compute gradient w/ autodiff

Autoencoder:

• bottleneck hidden layer w/ few nodes

• goal is to reconstruct input

• with single hidden layer & σ_1 linear equivalent to PCA

$$\hat{y} = \underline{w}_2 \underline{w}_1 y = \underline{w}_2 \underline{w}_1 \underline{U} \underline{S} \underline{V}^\top \leftarrow \text{SVD of } y$$

hidden layer
size K

to minimize

square loss,

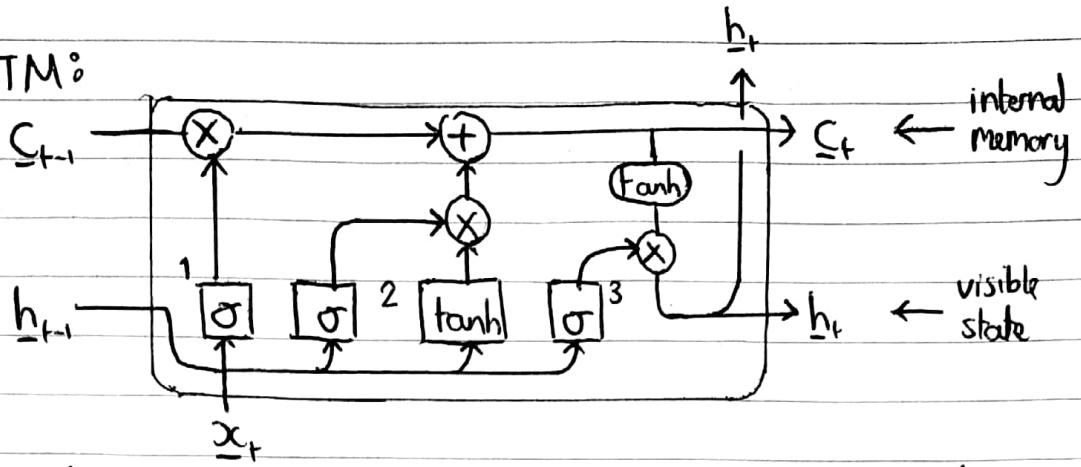
take largest

K singular values

or \underline{u}

$$\Rightarrow \underline{w}_1 \underline{w}_2 = \underline{U}_K \underline{U}_K^\top, \quad \hat{y} = \underline{U}_K \underline{U}_K^\top y, \quad \text{Same as PCA}$$

LSTM:



1: forget gate, 2: input gate, 3: output gate

- Special architecture for RNN, to help deal with exploding/vanishing gradients & long term memory

- train w/ backpropagation through time

Initialisation:

- aim to keep error $O(1)$

- scale layer inputs to mean 0, variance 1 ($x_i = \frac{x_i - \mu_i}{\sigma_i}$)

- sample weights w/ mean 0, variance = $1/D$ \leftarrow size of input \geq
- keep output mean 0, 1 variance

Decorrelating:

- makes opt Hebbian approx diagonal, can ease optimisation

$$\underline{x}^n = \underline{S}^{-1} \underline{U}^T \underline{x}^n, \text{ where } \underline{X} \underline{X}^T = \underline{U} \underline{S}^2 \underline{U}^T \text{ (so SVD of } \underline{X})$$

Auto Diff^g

given $g(\underline{x})$, find $\frac{\partial g}{\partial x_i} |_{\underline{x}}$

Forward:

- Complex arithmetic: approximate

$$g'(\underline{x}) = \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \text{Im}(g(\underline{x} + i\epsilon))$$

- Dual arithmetic: exact

$$g'(\underline{x}) = \text{Dual Part}(g(\underline{x} + \epsilon))$$

- both numerically inefficient & need to overload function to work w/ dual/complex arithmetic

Backwards:

- Find reverse schedule of nodes

- Starting from first node N (in reverse), $t_N = 1$

- Recurse back, $t_n = \sum_{c \in \text{ch}(n)} \frac{\partial g_c}{\partial g_n} t_c$ ($\text{ch}(n) = \text{children of } n$)

- t at root node giving total derivative wrt these variables

- efficient; $\leq S \times$ complexity of $g(\underline{x})$!

Dimension Reduction:

• PCA:

- minimize $E(\underline{\underline{B}}) = \sum_{n,d} (x_d^n - \sum_{j=1}^m y_j^n b_j^d)^2$

$$\frac{\partial E}{\partial y_k^n} = 0 \Rightarrow y_k^n = \sum_d x_d^n b_d \Rightarrow \underline{\underline{y}} = \underline{\underline{B}}^T \underline{\underline{X}}$$

$$E(\underline{\underline{B}}) = \text{tr}((\underline{\underline{X}} - \underline{\underline{B}} \underline{\underline{y}})^T (\underline{\underline{X}} - \underline{\underline{B}} \underline{\underline{y}})), \quad \underline{\underline{B}}^T \underline{\underline{B}} = \underline{\underline{I}}$$

$$= \text{tr}(\underline{\underline{X}}^T (\underline{\underline{I}} - \underline{\underline{B}} \underline{\underline{B}}^T)^2 \underline{\underline{X}}) = (\underline{\underline{X}} \underline{\underline{X}}^T (\underline{\underline{I}} - \underline{\underline{B}} \underline{\underline{B}}^T))$$

$$E(\underline{\underline{B}}) = \text{tr}(\underline{\underline{S}}) - \text{tr}(\underline{\underline{S}} \underline{\underline{B}} \underline{\underline{B}}^T) (N-1)$$

- minimize $-\text{tr}(\underline{\underline{S}} \underline{\underline{B}} \underline{\underline{B}}^T)$, enforce $\underline{\underline{B}}^T \underline{\underline{B}} = \underline{\underline{I}}$ w/ Lagrange Multiplier

$$\Rightarrow \underline{\underline{S}} \underline{\underline{B}} = \underline{\underline{B}} \underline{\underline{L}}$$

- enforce $\underline{\underline{L}}$ diagonal \Rightarrow cols of $\underline{\underline{B}}$ are eigenvectors of $\underline{\underline{S}} = \frac{1}{N-1} \underline{\underline{X}} \underline{\underline{X}}^T$

- eigendecomposition of $\underline{\underline{S}}$ $O(D^3)$, faster to do:

$$\underline{\underline{X}} = \underline{\underline{U}} \underline{\underline{D}} \underline{\underline{V}}^T, \quad \underline{\underline{B}} = \underline{\underline{U}} \quad (\text{SVD is } O(DN^2), O(D^2N))$$

- for sparse data, can iterate

$$\underline{\underline{x}} = \underline{\underline{S}} \underline{\underline{x}}_c, \text{ converges at principle eigenvector}$$

- for missing data,

$$E(\underline{\underline{B}}) = \sum_{n,d} y_d^n (x_d^n - \sum_{j=1}^m y_j^n b_j^d)^2$$

(mean of class)

- Fisher's Linear Discriminant

$$y^n = \underline{\underline{w}}^T \underline{\underline{x}}^n, \text{ project to minimize class overlap}$$

$$y_c \sim N(\underline{\underline{\mu}}_c, \underline{\underline{\sigma}}_c^2) \Rightarrow F(\underline{\underline{w}}) = \frac{(\underline{\underline{\mu}}_1 - \underline{\underline{\mu}}_2)^2}{\pi_1 \underline{\underline{\sigma}}_1^2 + \pi_2 \underline{\underline{\sigma}}_2^2} \quad \begin{aligned} \underline{\underline{\mu}}_c &= \underline{\underline{w}}^T \underline{\underline{M}}_c \\ \underline{\underline{\sigma}}_c^2 &= \underline{\underline{w}}^T \underline{\underline{S}}_c \underline{\underline{w}} \\ \pi_c &= N_c / N \end{aligned}$$

$$F(\underline{\underline{w}}) = \frac{\underline{\underline{w}}^T \underline{\underline{A}} \underline{\underline{w}}}{\underline{\underline{w}}^T \underline{\underline{B}} \underline{\underline{w}}} \quad \begin{aligned} \underline{\underline{A}} &= (\underline{\underline{\mu}}_1 - \underline{\underline{\mu}}_2)^T (\underline{\underline{\mu}}_1 - \underline{\underline{\mu}}_2) \\ \underline{\underline{B}} &= \pi_1 \underline{\underline{S}}_1 + \pi_2 \underline{\underline{S}}_2 \end{aligned}$$

$$\Rightarrow \frac{\partial F}{\partial \underline{\underline{w}}} = 0 \Rightarrow (\underline{\underline{w}}^T \underline{\underline{B}} \underline{\underline{w}}) \underline{\underline{A}} \underline{\underline{w}} = (\underline{\underline{w}}^T \underline{\underline{A}} \underline{\underline{w}}) \underline{\underline{B}} \underline{\underline{w}}$$

$$\Rightarrow \underline{\underline{w}} \propto \underline{\underline{B}}^{-1} (\underline{\underline{\mu}}_1 - \underline{\underline{\mu}}_2)$$

Clustering

- Can use k-means:

- hard clustering, scales poorly, requires spherical clusters w/ equal density

- Mixture models:

- $$p(\underline{x}) = \sum_c p(\underline{x}|c) \cdot p(c)$$
 (PSD, symmetric)

- $$\text{eg GMM: } \underline{x}|c \sim N(\underline{\mu}_c, \underline{\Sigma}_c^2)$$

- learn $\underline{\mu}_c, \underline{\Sigma}_c^2$ by maximising likelihood (EM algorithm)

$$L = \sum_n \log \sum_c p(c) \cdot p(\underline{x}|c) \quad \left(p(\underline{x}|c) = \frac{1}{\sqrt{\det(2\pi\underline{\Sigma}_c)}} \exp \left(-\frac{1}{2} \frac{\underline{x}^n - \underline{\mu}_c}{\underline{\Sigma}_c} \right) \right)$$

- need to stop $\underline{\Sigma}_c$ being too small, else $p(\underline{x}|c)$ can become infinite for $\underline{x} = \underline{\mu}_c$

- Gaussian sensitive to outliers; use t-distributions

- Spectral clustering:

- data as markov chain

$$p(i|j) = \frac{A_{ji}}{\sum_k A_{jk}} \quad (A_{ji} \text{ is similarity of } j \text{ & } i)$$

- find equilibrium dist & cluster with this (k-means)

Fast Nearest Neighbour

- Metric distance:

- $d(x, y) \geq 0$

- $d(x, y) = d(y, x)$

- $d(x, y) = 0 \Rightarrow y = x$

- $d(x, y) \leq d(x, z) + d(z, y)$

neighbours

- Orchard:

- precompute distance matrix & ordered list for each \underline{x}^i

- for query q^8

- pick initial \underline{x}^i , step along its list to \underline{x}^i

- if $d(\underline{x}^i, q) < d(\underline{x}^i, q)$, go to \underline{x}^i 's list

- if $d(\underline{x}^i, q) \geq \frac{1}{2} d(\underline{x}^i, q)$, stop early, return \underline{x}^i

- AESA:

- precompute distance matrix

- $\mathcal{X} = \text{set of } d(q, \underline{x}^i) \text{ calculated already}$

(lower bound
on $d(q, \underline{x}^i)$)

- pre-eliminate \underline{x}^i if: ~~if~~

$$\min_{i \in \mathcal{X}} \{d(q, \underline{x}^i)\} \leq \max_{i \in \mathcal{X}} \{d(q, \underline{x}^i) - d(\underline{x}^i, \underline{x}^i)\}$$

- Buoys:

- small set of points $\underline{b}^1, \dots, \underline{b}^m$

- compute distance matrix only for buoys

- for \underline{x}^j :

$$U_j = \min_m \{d(q, \underline{b}^m) + d(\underline{b}^m, \underline{x}^j)\}$$

$$L_j = \max_m \{d(q, \underline{b}^m) - d(\underline{b}^m, \underline{x}^j)\}$$

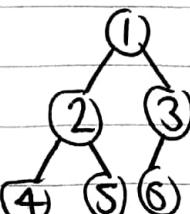
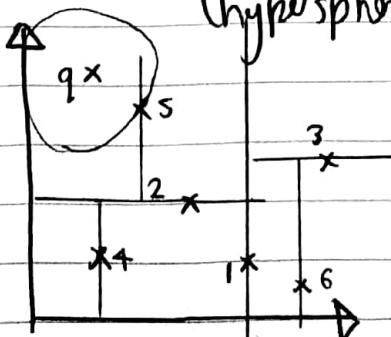
- pre-eliminate all points with $U_j > \min_i \{L_i\}$

- then feed to AESA / Orchard

- k-d tree:

- recursively divide space, alternating dimensions

- query goes up, checking if need to check siblings
(hypersphere around sibling ~~already seen~~ query)



$q \rightarrow$ check 5
 skip 4
 check 2
 skip 3
 check 1 \Rightarrow return 5

COMP9009 - Applied ML

1a: Intro to Supervised Learning

- Loss: cost of making an error ($U = -L$)
- For input x^* , we learn predictive distribution $P(c|x^*)$ & function, $c(x^*)$
- Expected utility = $U(c(x^*)) = \sum_{c^{\text{true}}} U(c^{\text{true}}, c(x^*)) \cdot P(c^{\text{true}}|x^*)$

try to maximise this

- How to define $P(c|x^*)$ & $c(x^*)$?

- Empirical Risk: empirical dist $P(c|x^*)$, Parametric $c(x^*|\theta)$
- Bayesian Decision: non-trivial $P(c|x^*)$, unconstrained $c(x^*)$

- Utility functions:

• Utility General utility matrix:

$$U_{i,j} = U(c^{\text{true}}=i, c^{\text{pred}}=j) \quad \begin{matrix} \text{can be V. Asymmetric;} \\ \text{eg cancer diagnosis} \end{matrix}$$

$$U(c^*=j) = \sum_i U_{i,j} \cdot P(c^{\text{true}}=i|x^*)$$

• Squared Loss:

$$L(y^{\text{true}}, y^{\text{pred}}) = (y^{\text{true}} - y^{\text{pred}})^2$$

$$L(y^{\text{pred}}) = \int (y^{\text{true}} - y^{\text{pred}})^2 p(y^{\text{true}}|x) dy^{\text{true}}$$

$$(\text{maximised by } y^{\text{pred}} = \int y^{\text{true}} p(y^{\text{true}}|x) dy^{\text{true}})$$

- Choice of loss can influence ease of training; logistic regression model for classification w/ log loss is convex, squared loss is non-convex

• Empirical Risk Approach:

$$D = \{(x^n, c^n), n=1 \dots N\}$$

• approximate $P(c|x)$ by empirical dist, $P(c, x|D) = \frac{1}{N} \sum_n \delta(c, c_n) \delta(x, x_n)$

$$\Rightarrow \langle U(c, c(x)) \rangle_{P(c, x|D)} = \frac{1}{N} \sum_n U(c_n, c(x_n)) \leftarrow \begin{matrix} \text{Empirical} \\ \text{expected} \\ \text{utility} \end{matrix}$$

• For x_n^* in train set, $c(x_n) = c_n$, but ~~this x~~ ^{outside train set} have $c(x)$ undefined!

• For x outside train set:

use parametric function $c(x|\theta)$

$$\text{eg 2 class problem, } c(x|\theta) = \begin{cases} 1 & \text{if } \vec{\theta}^T \vec{x} + \theta_0 \geq 0 \\ 2 & \text{o/w} \end{cases}$$

$$R(\theta|D) = \frac{1}{N} \sum_n L(c_n, c(x_n|\theta))$$

$$\theta_{\text{opt}} = \underset{\theta}{\operatorname{arg\,min}} R(\theta|D)$$

• Overfitting:

• too flexible $c(x|\theta)$ overfits training data

• add penalty term, $P(\theta)$ (eg $P(\theta) = \|\vec{\theta}\|^2$)

$$\text{minimise } R(\theta|D) = R(\theta|D) + \lambda P(\theta)$$

λ set by validation

• Cross-validation:

$$D = \boxed{\quad} \boxed{\quad} \boxed{\quad} \dots \dots \boxed{\quad} \text{ Sp}$$

• split data into k folds

• train on $k-1$ folds, test on held out fold

• repeat k times & average loss/accuracy ~~loss~~

• Pros & Cons:

+ large training data, empirical dist tends towards true

+ discriminant function chosen to minimize risk

- poor for smaller amount of train data

- need to train for a different loss function

- no inherent confidence in prediction

- Bayesian Decision Approach:

- fit model $P(c, x | \Theta)$ to train data (e.g. maximum likelihood)
 - could overfit, maybe use regularization
- determine $c(x)$ given by maximizing expected utility

$$U(c(x^*)) = \sum_c U(c, c(x^*)) P(c | x, \Theta)$$

$$c(x^*) = \underset{c(x^*)}{\operatorname{argmax}} U(c(x^*))$$

- Discriminative-Generative Approach:

$$p(c, x | \Theta) = p(x | c, \Theta_{x|c}) \cdot p(c | \Theta_c) \quad (\underset{x}{\text{class generates}})$$

- + easy to specify prior in $p(x | c)$
- doesn't directly target classification model, $p(c | x)$;
can be non-trivial to find $p(c | x)$ now if
data is complex

- Discriminative Approach:

$$p(c, x | \Theta) = p(c | x, \Theta_{c|x}) \cdot p(x | \Theta_x)$$

- + directly addresses $p(c | x)$, modelling decision boundary (which may be far easier to model than distribution of data within a class)
- hard to include prior knowledge

- Pros & Cons:

- + conceptually very good, trying to model environment independent of decision process
- + if $p(x, c | \Theta)$ correct, optimal
- if $p(x, c | \Theta)$ poor, prediction could be highly inaccurate
 - often use regularization for $p(x, c | \Theta)$ to help avoid this

1b: Maths Refresher

Linear Algebra

$$A = \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \dots & a_{mn} \end{pmatrix}$$

$$[A + B]_{ij} = a_{ij} + b_{ij}$$

$$[AB]_{jk} = \sum_{i=1}^l a_{ij} b_{jk} \quad (\text{if } AB = BA, A \& B \text{ commute})$$

$$I_{ij} = \delta_{ij}, \quad \delta_{ij} = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{otherwise} \end{cases} \quad I = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

$$[B^T]_{kj} = b_{jk} \quad (ABC)^T = C^T B^T A^T$$

Vectors: can be $n \times 1$ column matrix $\vec{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$

$$\vec{w} \cdot \vec{x} = \sum w_i x_i = \vec{w}^T \vec{x}$$

$$|\vec{x}|^2 = \vec{x}^T \vec{x} = \sum x_i^2$$

$$\det(A^T) = \det(A)$$

$$\det(AB) = \det(A) \det(B)$$

$$\det(I) = 1$$

$$\det(A^{-1}) = \frac{1}{\det(A)}$$

~~Linear~~ $\rightarrow A^{-1}A = I = AA^{-1}$ $(\text{if } A^{-1} \text{ doesn't exist, } A \text{ is singular})$
 $(AB)^{-1} = B^{-1}A^{-1}$

~~Linear~~ system $\rightarrow A\vec{x} = \vec{b} \Rightarrow \vec{x} = A^{-1}\vec{b}$

$X = [\vec{x}_1, \dots, \vec{x}_n]$, Rank is max # linearly independent cols or rows

$\vec{x}_1, \dots, \vec{x}_m$ linearly independent if only solution to $\sum_i \alpha_i \vec{x}_i = \vec{0}$ is $\alpha_i = 0$ for all i

$$A^T A = I = A A^T \Rightarrow A^T = A^{-1}, A \text{ is orthogonal}$$

eigenvectors
↓ eigenvalues

$A\vec{e} = \lambda\vec{e}$, for $n \times n$ matrix, n eigenvalues & vectors
 $(A - \lambda I)\vec{e} = \vec{0}$
 $\Rightarrow \lambda$ is eigenvalue if $\det(A - \lambda I) = 0$

$$A = \sum_{i=1}^n \lambda_i \vec{e}_i \vec{e}_i^T$$

eigenvectors orthogonal to each other

$$A = E \Lambda E^{-1} \quad (E \text{ matrix of eigenvectors, } \Lambda = \text{diag}(\lambda_i))$$

Singular
Value
decomposition $\rightarrow X = USV^T$

$$U: n \times n, \quad U^T U = I_n$$

$$V: p \times p, \quad V^T V = I_p$$

$$X: n \times p$$

$S: n \times p$ & diagonal,
all +ve, &
ordered (largest top-left)

if $\forall \vec{x}, \vec{x}^T A \vec{x} > 0 \quad (\vec{x} \neq \vec{0}) \Rightarrow A$ positive definite

$$\text{trace}(A) = \sum a_{ii} = \sum \lambda_i \quad \forall \log \det(A)$$

$$\det(A) = \prod \lambda_i$$

$$\text{trace}(\log A) = \log \det(A)$$

1b: Maths Refresher

Calculus:

$$\frac{df}{dx} = \lim_{\delta \rightarrow 0} \frac{f(x+\delta) - f(x)}{\delta} = f'(x) \quad \begin{aligned} & \left(\frac{df}{dx} \approx \frac{f(x+\delta) - f(x)}{\delta} + O(\delta^2) \right) \\ & \approx \frac{f(x+\delta) - f(x-\delta)}{2\delta} + O(\delta^3) \end{aligned}$$

Taylor series

$$\rightarrow f(x) = f(0) + x \frac{df}{dx} + \frac{x^2}{2} \frac{d^2 f}{dx^2} + \dots$$

chain rule

$$\rightarrow \frac{df(g(x))}{dx} = \frac{df(y)}{dy} \Big|_{y=g(x)} \quad \frac{dg}{dx} = \frac{df}{dx} \frac{dg}{dx}$$

$$\frac{d}{dx}(f+g) = \frac{df}{dx} + \frac{dg}{dx}, \quad \frac{d}{dx}(fg) = f \frac{dg}{dx} + g \frac{df}{dx}$$

$$\frac{df}{dx} = \frac{\partial f}{\partial x} + \frac{\partial f}{\partial y} \frac{\partial y}{\partial x} \quad \begin{array}{c} x \downarrow \frac{df}{dx} \\ \downarrow \frac{\partial f}{\partial y} \\ g \end{array}, \quad \text{where } f(x) \text{ includes } g(x)$$

Jacobian $\rightarrow \nabla f(\vec{x}) = \vec{g}(\vec{x}) = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix}$, points along direction of max change

Hessian $\rightarrow H_f(\vec{x}) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{pmatrix}$

$$\frac{\partial}{\partial A} \text{trace}(AB) = B^T$$

$$\partial \log \det A = \partial \text{trace}(\log A) = \text{trace}(A^{-1} \partial A)$$

$$\frac{\partial}{\partial A} \log \det A = A^{-T}$$

$$\partial A^{-1} = -A^{-T} \partial A A^{-1}$$

numerical approximation



- Convex analysis:

- $f(\vec{x})$ convex if, for any \vec{x}, \vec{y} , $0 < \lambda < 1$,

$$f(\lambda \vec{x} + (1-\lambda) \vec{y}) \leq \lambda f(\vec{x}) + (1-\lambda) f(\vec{y})$$

(or if Hessian is positive definite for all points)

- only have 1 minimum & that is only global point

- f, g convex $\Rightarrow f+g, \cancel{f \circ g}, f(Ax+b)$ convex
 $g(g(x))$ convex if g increasing function

- Numerical issues:

- mathematical identities may not hold on computer; need to check & fix at key points in code

- e.g. $S = \exp(a) + \exp(b)$ may well overflow, leaving $\log(S)$ undefined

- use logsumexp, ~~approx~~

$$\log(S) = a + \log(1 + \exp(b-a)) \quad \begin{cases} \text{if } a \geq b, \text{ in general} \\ \text{take the max} \end{cases}$$

$$\text{logsumexp}(\vec{x}) = m + \log\left(\sum \exp(x_i - m)\right), \quad m = \max(\vec{x})$$

- Distributions:

- Multi-variate Gaussian

$$P(\vec{x} | \vec{\mu}, \Sigma) = \frac{1}{\sqrt{\det(2\pi\Sigma)}} \exp\left[-\frac{1}{2}(\vec{x} - \vec{\mu})^\top \Sigma^{-1} (\vec{x} - \vec{\mu})\right]$$

- Σ is covariance, & symmetric

$$\Rightarrow \Sigma = E \Lambda E^\top \quad (\text{eigenvalue decomposition})$$

- can view multi-variate Gaussian as shifted, scaled & rotated standard Gaussian

Decomposing \rightarrow $\vec{y} = \Sigma_x^{-1/2} (\vec{x} - \vec{\mu}_x) \sim N(\vec{0}, I)$

2: Regression

- $S = \{(\vec{x}_i, y_i), i=1, \dots, N\}$ sampled iid from $P(\vec{x}, y)$
want to learn mapping from \vec{x} to y

- Least Squares:

$$E(\vec{w}) = \sum_i (y_i - \vec{w}^T \vec{\phi}_i), \quad \vec{\phi}_i = (1, \vec{x}_i)^T = \vec{\phi}(\vec{x}_i)$$

\hat{w} , minimiser of $E(\vec{w})$, satisfies $(X^T X) \hat{w} = X^T \vec{y}$

$$X = \begin{pmatrix} \vec{x}_1^T \\ \vdots \\ \vec{x}_n^T \end{pmatrix}, \quad \vec{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

- can fit polynomial w/ different basis functions, eg
for cubic $\vec{\phi}(x) = (1, x, x^2, x^3)^T$

- Regularisation:

- Should first scale input dimensions to unit variance

$$x_{id} \leftarrow \frac{x_{id}}{\sigma_d}, \quad \sigma_d^2 = \frac{1}{N} \sum_i (x_{id} - \mu_d)^2, \quad \mu_d = \frac{1}{N} \sum_i x_{id}$$

- now minimise $E'(\vec{w}) = E(\vec{w}) + \lambda R(\vec{w})$

- L_2 -Regularisation:

$$R(\vec{w}) = |\vec{w}|^2 = \vec{w}^T \vec{w} = \sum w_i^2$$

$$\hat{w} = (\sum_i \vec{\phi}_i (\vec{\phi}_i)^T + \lambda I)^{-1} \sum_i y_i \vec{\phi}_i$$

- + easy analytic solution

- doesn't encourage sparse solutions; prefers many small coefficients (eg collinear x , keeps both)

- L_1 -Regularisation

$$R(\vec{w}) = |\vec{w}| = \sum w_i$$

- + encourages sparse solutions

only at origin, can just project & not worry about sometimes

- Non-differentiable; need special optimisation routines

- Doesn't deal w/ co-linear x , but not usually an issue

- Auto-Regressive models:

- time series of V_1, \dots, V_{t-1}

- predict V_t using V_{t-1}, \dots, V_{t-L}

$$V_t \approx \sum_{l=1}^L a_l V_{t-l} = \vec{a}^T \vec{\Phi}_t, \quad \vec{\Phi}_t = (V_{t-1}, \dots, V_{t-L})^T$$

- can train a_l before, w/ LS & L_1/L_2 -regularisation

- Probabilistic View:

$$V_t = \vec{a}^T \vec{\Phi}_t + \eta_t, \quad \eta_t \sim N(\eta_t | \mu, \sigma^2) \quad \text{"innovation level"}$$

- L^{th} order Markov model

$$P(V_{1:T}) = \prod_{t=1}^T P(V_t | V_{t-1}, \dots, V_{t-L})$$

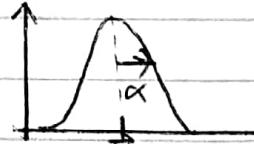
$$P(V_t | V_{t-1}, \dots, V_{t-L}) \sim N(V_t | \vec{a}^T \vec{\Phi}_t, \sigma^2)$$

- $L(\vec{a}) = \sum_{t=1}^T P(V_t | \vec{\Phi}_t)$

Maximised by $\vec{a} = (\sum_t \vec{\Phi}_t \vec{\Phi}_t^T)^{-1} \sum_t V_t \vec{\Phi}_t$

- Radial basis functions:

- eg $\phi_i(\vec{x}) = \exp\left(-\frac{1}{2\alpha^2} (\vec{x} - \vec{m}_i)^2\right)$



- adding together these, with weights, allows function fitting

- curse of dimensionality:

- if in 1D need 16 ^{RBF} functions to work well, need 16^n in n -dimensions for same performance

- Solution: let centres move (similar to neural networks)

- or use Gaussian processes

3: Classification

- $S = \{(\vec{x}_i, y_i), i=1, \dots, n\}$, y_i discrete class labels
- Generative models:
 - k-NN:
 - to classify \vec{x}^* , look at k nearest points in train set & do majority vote
 - distance
 - squared Euclidean: $d(\vec{x}, \vec{x}') = (\vec{x} - \vec{x}')^T (\vec{x} - \vec{x}')$
 - if length scales of diff components of \vec{x} vary largely, will be dominated by one, losing information
 - Mahalanobis distance: $d(\vec{x}, \vec{x}') = (\vec{x} - \vec{x}')^T \Sigma^{-1} (\vec{x} - \vec{x}')$
 - + rescales components of \vec{x}
 - + simple, smooth decision boundary
 - need to store whole training set, distance calculation slow, unclear how to deal w/ missing data or incorporate priors
 - Probabilistic View:
 - each train point is centre of Gaussian w/ width σ^2
 - $\sigma^2 \rightarrow 0 \Rightarrow$ only nearest neighbour has effect, $\sigma^2 \rightarrow \infty \Rightarrow$ others have effect
 - training point classified by looking at these Gaussian generated by training points
 - giving even smoother k-NN
 - if neighbours far, far neighbours classified very confidently
 - Naive Bayes: $P(c|\vec{x}^*) = \frac{P(\vec{x}^*|c)P(c)}{P(\vec{x}^*)} = \frac{P(\vec{x}^*|c)P(c)}{\sum P(\vec{x}^*|c)P(c)}$
 - + very fast to train; can just count for discrete attributes. Deals with missing data easily
 - overly confident with low counts (solution: add pseudo counts for each class)

• Discriminative Models:

- parametric classification: classify \vec{x} with $c(\vec{x} | \vec{\theta})$, learn $\vec{\theta}$ from data
- set $\vec{\theta}$ to minimize loss on train set, eg $L(c, \hat{c}) = \begin{cases} 0 & \text{if } c = \hat{c} \\ 1 & \text{otherwise} \end{cases}$ (Zero-one loss) surface, hard to train

logistic
sigmoid

• Logistic Regression:

$$p(c=1 | \vec{x}) = \sigma(b + \vec{x}^T \vec{w}), \quad p(c=0 | \vec{x}) = 1 - p(c=1 | \vec{x})$$

↑
for multiclass,
change to softmax

$$g(x) = \frac{e^x}{\sum e^x}$$

• Decision boundary:

$$\vec{x} \text{ where } p(c=1 | \vec{x}) = p(c=0 | \vec{x}) \Rightarrow b + \vec{x}^T \vec{w} = 0$$

• linearly separable if all of each class lies on either side

of decision boundary

log likelihood
(assume data iid)

• if not linearly separable, use of basis function can help

$$\rightarrow L(\vec{w}, b) = \sum_i c_i \log \sigma(b + \vec{w}^T \vec{x}_i) + (1 - c_i) \log(1 - \sigma(b + \vec{w}^T \vec{x}_i))$$

• optimize by gradient ascent or Newton (faster)

• convex, so will reach global optimum

• if linearly separable, weights will forever grow; need to stop early or add penalty term

• Support Vector Machines:

$$\text{• classes } \in \{-1, 1\}, \quad c(\vec{x}) = \begin{cases} 1 & \text{if } \vec{w}^T \vec{x} + b \geq 0 \\ -1 & \text{if } \vec{w}^T \vec{x} + b < 0 \end{cases}$$

• goal: minimize $\frac{1}{2} \vec{w}^T \vec{w}$

subject to $y_i (\vec{w}^T \vec{x}_i + b) \geq 1, \quad i=1, \dots, n$

• if data not linearly separable:

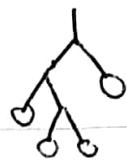
$$\text{Minimise } \frac{1}{2} \vec{w}^T \vec{w} + \frac{1}{2} C \sum_i \epsilon_i^2$$

subject to $y_i (\vec{w}^T \vec{x}_i + b) \geq 1 - \epsilon_i, \quad i=1, \dots, n$

• Kernels:

• if use basis function $\phi(\vec{x})$, only need to be able to calculate $\langle \phi(\vec{x}_1), \phi(\vec{x}_2) \rangle$; can skip the basis function and just use a kernel; ie an implicit basis function

• effectively performs SVM in higher (or infinite) dimensional space



- Decision trees:

- learn series of decisions, such that data at each node as 'pure' as possible
 - + fast to train & interpretable
 - unstable & limited

- Ensemble:

- combine predictions of several classifiers

- Bagging:

- generate B datasets by sampling (with replacement) from original

- train B models on these, & average or vote for predictions

- Random forest is roughly bagging applied to decision trees

- Summary:

- Discriminative:

- K-NN + simple, intuitive

- slow ($O(n)$), need good distance function, struggle w/ missing data

- Naive Bayes + fast to train & apply, easily deals w/ missing data

- weak

- Discriminative:

- Logistic Regression + simple & fast to train, fast to apply, ($O(p)$), concave optimisation, can use kernels, scales well

- struggles w/ missing data

- SVM + good generalisation, fast to apply, concave

- doesn't scale well, struggles w/ missing data

- Decision Tree + fast to train, interpretable, - very weak, struggles w/ missing data

- Random Forest + fast to train, good performance

- less interpretable, struggle w/ missing data

4: Optimisation

- Gradient Descent & First order methods:

$$\vec{x}_{k+1} = \vec{x}_k + \epsilon \nabla f(\vec{x}_k) \quad (\text{decrease } f \text{ as } f(\vec{x}_{k+1}) \approx f(\vec{x}_k) - \epsilon \|\nabla f(\vec{x}_k)\|^2)$$

- can get stuck in local minima, & learning rate too big
may diverge

- Convergence:

- assume f convex, optimum is finite, f twice differentiable
- $f(\vec{x}_T) - f(\vec{x}^*) \leq \frac{1}{2\epsilon T} (\vec{x}_1 - \vec{x}^*)^2$

\Rightarrow error falls as $1/T$ ← iterations of grad descent

- Momentum:

- reduce zig-zag behaviour by updating \vec{x} w/ moving average of gradients

$$\begin{aligned}\tilde{g}_{k+1} &= \mu_k \tilde{g}_k - \epsilon \nabla f(\vec{x}_k) \\ \vec{x}_{k+1} &= \vec{x}_k + \tilde{g}_{k+1}\end{aligned}$$

ie saddles

- Converges faster, helps carry us across flat regions too

- very good for noisy gradients

- Nesterov:

- Same as momentum, except

$$\tilde{g}_{k+1} = \mu_k \tilde{g}_k - \epsilon \nabla f(\vec{x}_k + \mu_k \tilde{g}_k)$$

effectively, do momentum, then think 'how can I get an even better update using gradient at current point'

- ie use gradient of point we will move to, rather than current point

$$\mu_k = 1 - 3/(k+5), \text{ convergence } f(\vec{x}_k) - f^* \leq \frac{\epsilon}{k^2} \text{ (comp. w/ } \frac{1}{k})$$

- very simple & easy to implement

- need suitably small ϵ to guarantee convergence

- coordinate system dependant; transformation of coordinates leads to different path

↑
eg changing units

• Line Search & Conjugate Gradients:

• Line Search:

- pick direction, find optimum in this direction, repeat
- if quadratic function, $\vec{d} = -\nabla f(\vec{x}_0)$ zig-zags; better to line-search independently along each axis

• Conjugate Gradients:

- perform line search updates in conjugate directions sequentially, building up set of conjugate vectors
- after d iterations, will have found overall optimum (to a quadratic function)

$$\cdot k=1; \vec{x}_1 = \vec{x}_0; \vec{p}_1 = -\vec{g}_1$$

while $\vec{g}_k \neq 0$:

$$\alpha_k = \arg \min_{\alpha} f(\vec{x}_k + \alpha \vec{p}_k) \leftarrow \begin{array}{l} \text{line} \\ \text{search} \end{array}$$

$$\vec{x}_{k+1} = \vec{x}_k + \alpha_k \vec{p}_k \leftarrow \text{update } x$$

$$\beta_k = \vec{g}_{k+1}^\top \vec{g}_{k+1} / (\vec{g}_k^\top \vec{g}_k) \quad] \text{ find new}$$

$$\vec{p}_{k+1} = -\vec{g}_{k+1} + \beta_k \vec{p}_k \quad] \text{ conjugate direction}$$

$\vec{p}_{k+1} \text{ conjugate}$
to $\vec{p}_1, \dots, \vec{p}_k$

• Higher order methods:

• Newton's method:

Taylor expansion to 2nd order $\longrightarrow f(x + \Delta) = f(x) + \Delta^\top \nabla f + \frac{1}{2} \Delta^\top H_f \Delta + O(|\Delta|^3)$

$$\Rightarrow x_{k+1} = x_k - \varepsilon H_f^{-1} \nabla f \quad (\text{converge in 1 step for quadratic functions, } \varepsilon=1)$$

+ invariant under coord transform

- not storing Hessian & calculating $H^{-1} \nabla f$ expensive

- not guaranteed to produce downhill step; can do line search in $H^{-1} \nabla f$ direction instead

• Gauss-Newton:

for objective function $f(\vec{x}) = L(\vec{h}(\vec{x}))$,

approaches Newton as we approach optimum, as G approaches H

$$G_{ij} = \sum_{k=1}^m \frac{\partial h_k}{\partial x_i} \frac{\partial^2 L}{\partial h_k \partial h_i} \frac{\partial h_k}{\partial x_j} \quad \begin{array}{l} \text{can approximate } H, \\ \text{is positive semidefinite} \end{array}$$

$$\Rightarrow \text{update } \vec{x}_{t+1} = \vec{x}_t - \varepsilon G^{-1} \vec{g}, \text{ with } g_i = \frac{\partial L}{\partial x_i}$$

- Levenberg - Marquardt algorithm:

- Gauss-Newton needs to solve $C\Delta = \nabla E$
- C can be rank deficient or ill conditioned
- regularize, & instead solve $(C + \lambda I)\Delta^* = \nabla E$

Gauss-Newton
for KL loss

- Natural Gradient:

$$E(\Theta) = KL(p(x) \| q(x|\Theta)) = -\langle \log q(x|\Theta) \rangle + \text{const}$$

gradient $\rightarrow g_i = -\left\langle \frac{\partial}{\partial \Theta_i} \log q(x|\Theta) \right\rangle, C_{ij} = \left\langle \frac{\partial}{\partial \Theta_i} \log q(x|\Theta) \cdot \frac{\partial}{\partial \Theta_j} \log q(x|\Theta) \right\rangle$

update $\rightarrow \Theta(t+1) = \Theta(t) - \varepsilon C^{-1} g$

curvature matrix

+ invariant under parameter transforms

Sa: Large Scale Learning

• Linear Models:

- least squares is $O(\dim(\vec{\theta})D^3)$ time to solve, & requires $O(D^2)$ space
- ~~can~~ don't need much space to store ~~gradient~~ ^{comput} gradient

$$E(\vec{\theta}) = \sum (y_n - \vec{\theta}^T \vec{x}_n)^2 + \lambda \vec{\theta}^T \vec{\theta} \quad \text{(where } d \text{ is density factor of sparse } \vec{x}, \text{ i.e. each } \vec{x} \text{ has max } dD \text{ non-zero elements)}$$
$$\frac{\partial E}{\partial \theta_i} = -2 \sum (y_n - \vec{\theta}^T \vec{x}_n) x_{ni} + 2\lambda \theta_i$$

• Online stochastic gradient descent:

- scale objective $E(\vec{\theta}) = \frac{1}{N} \sum_n (y_n - \vec{\theta}^T \vec{x}_n)^2 + \gamma \vec{\theta}^T \vec{\theta}$

$$\frac{\partial E}{\partial \theta_i} = -2 \frac{1}{N} \sum_n (y_n - \vec{\theta}^T \vec{x}_n) x_{ni} + 2\gamma \theta_i$$

batch size

$$\downarrow \quad \quad \quad \approx -2 \frac{1}{N'} \sum_n (y_{n'} - \vec{\theta}^T \vec{x}_{n'}) x_{n'i} + 2\gamma \theta_i$$
$$N' < N$$

- extreme case $N' = 1$; treat inputs as a stream
- + robust to synchronisation issues in parallel implementations

• Higher order methods:

- use better search directions + line search to converge faster
- e.g. conjugate gradients

5b: Neural Net Introduction

- Vector input \vec{x} , $\vec{y} = g(\vec{x}|W) = \sigma_L(W_L \vec{h}_{L-1})$

- where $\vec{h}_1 = \sigma_1(W_1 \vec{h}_{1-1})$, $1=2,\dots,L-1$, $\vec{h}_1 = \sigma_1(W, \vec{x})$
- σ_1 is the transfer function (normally invertible & monotonic)
- common to include bias terms too (so $\vec{h}_1 = \sigma(W_1 \vec{h}_{1-1} + b_1)$)

$$\vec{x} \rightarrow \vec{h}_1 \rightarrow \dots \rightarrow \vec{h}_{L-1} \rightarrow \vec{y}$$

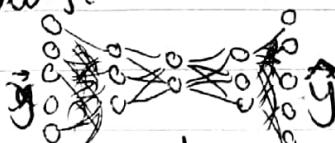
- if σ is linear, whole network is linear & additional layers play no role

• Regression:

- minimise loss function (eg squared loss)
- optimise using SGD or higher order methods
 - gradient easy to calculate w/ AutoDiff
- hard to initialise weights well; variety of approaches to including layerwise training
- common to include penalty (eg L_2 -regularisation) on weights
 - otherwise very easy to overfit

• Specific architectures:

• Autoencoder:



$$E(W) = \sum_n (\vec{y}_n - g(\vec{y}_n|W))^2$$

- try to reconstruct input
- bottleneck layer forces network to find low-dimensional representation of data
- equivalent to PCA when using network $y \rightarrow h \rightarrow \vec{y}$ w/ linear transfer function
- can only beat PCA with a non-linear function at output
- if $N = \text{size of dataset}$ units in bottom hidden layer can always perfectly reconstruct output even w/ bottleneck size = 1

- Convolutional NNs:
 - allows for translational (or higher order) of invariance; good for image processing
 - usually then max-pool after conv layer, to detect if any of neurons responded strongly

- NNs in NLP:
 - Word embeddings:
 - replace 1-hot euclidean embeddings w/ embeddings in lower dimensional space
 - can learn diff embeddings dependent on goal; eg next word prediction
 - natural geometry tends to appear, eg $v_{WOMAN} - v_{MAN} \approx v_{QUEEN} - v_{KING}$

- Time Series Models:
 - Recurrent nets:
 - hidden activation at time t (h_t) (& possibly y_t) depend on h_{t-1} , x_{t-1} & y_{t-1}
 - gradient computation:
$$E(A, B, C) = \sum_t (y_t - f(h_t; C))^2, \quad h_t = g(x_t, h_{t-1}; A, B)$$

- RTRL - Real time recurrent learning:
 - compute gradient w/ single forward pass in time ($\frac{\partial h_t}{\partial \theta} = f(x_t, h_{t-1})$)
 - $O(tH^3)$ space ($H = \text{size of hidden layer}$)
 \Rightarrow normally impractical

- BPTT - Backprop Through Time:
 - use parameter tying; calculate gradient treating each layer params as independent, then enforce evaluate expression by setting all to the same value
 - efficient & exact

- hard to train, usually use special architectures (eg LSTM)

5b: Neural Net Introduction

Generative Models:

- parametrise joint distributions (ie describes how points are generated)

$$P(x, y) = P(x|y) \cdot P(y)$$

- can then find information about y from $P(y|x)$

- Variational inference:

$$P(x, y|\Theta) = \int_y P(x|y, \Theta) \cdot P(y)$$

- We want to find Θ to maximise this

- use a variational distribution $q(y|x, \phi)$ to help estimate $P(y)$, & optimise wrt Θ & ϕ

Training nets:

- $E(\vec{w})$ is usually complex, not convex, many non-trivial local optima

- Solution quality depends lots on weight initialisation

- + Very powerful when trained well & right

7a: Automatic Differentiation

- take $g(\vec{x})$ & return exact value for $\frac{\partial g}{\partial x_i} \Big|_{\vec{x}}$
 - can do in less than $5x$ time to calculate $g(\vec{x})$
- Symbolic differentiation:
 - can give computationally inefficient results, e.g. need to repeat calculations unnecessarily for $\frac{\partial g}{\partial x_1} \Big|_{\vec{x}}$ & $\frac{\partial g}{\partial x_2} \Big|_{\vec{x}}$
- Forward differentiation:
 - Complex arithmetic:
$$\text{Dual exact } g(x+\epsilon) = \text{DualPart}(g(x+\epsilon))$$

$$\text{Complex approx } g'(x) = \frac{g(x+\epsilon) - g(x)}{\epsilon}$$
 - can do similar w/ complex arithmetic, keeping the imaginary part

$$(\text{e.g. } g(x) = x^2 \Rightarrow g(x+\epsilon) = x^2 + \boxed{2\epsilon x})$$

- holds for any smooth function, need to overload every function to work with dual arithmetic
- not numerically efficient

• Reverse differentiation:

1. Find reverse ancestral schedule of nodes
2. Start from 'first node' in reverse schedule, $t_n = 1$

3. For next node n in schedule:

$$t_n = \sum_{c \in \text{ch}(n)} \frac{\partial g_c}{\partial g_n} t_c, \text{ where } \text{ch}(n) \text{ is children of } n$$

4. Total derivatives wrt root nodes are given by values of t at these nodes

e.g.

$$g(x_1, x_2) = \cos(\sin(x_1, x_2))$$
$$x_1 \quad \quad \quad x_2$$
$$\frac{\partial g}{\partial x_1} = x_2 \quad \quad \quad \frac{\partial g}{\partial x_2} = x_1$$
$$f_1 \quad \quad \quad f_2$$
$$\downarrow \frac{\partial g_1}{\partial f_1} = \cos(f_1)$$
$$f_1 \quad \quad \quad f_2$$
$$\downarrow \frac{\partial g_2}{\partial f_2} = -\sin(f_2)$$
$$f_3$$
$$\uparrow t_{x_1} = x_2 t_3, \uparrow t_{x_2} = x_1 t_3$$
$$\uparrow t_3 = \cos(-\sin(f_2)) \cos(f_1) \cdot (-\sin(f_2))$$
$$\uparrow t_2 = -\sin(f_2)$$
$$\uparrow t_1 = 1$$

- Hessian-Vector Product:

- in Newton optimisation, update using $\vec{x} = H^{-1} \vec{g}$
- in practice, find update by solving $H\vec{x} = \vec{g}$
 \Rightarrow need $H\vec{x}$

derivative in

direction \vec{v} $\rightarrow D_{\vec{v}}(f) = \lim_{\delta \rightarrow 0} \frac{f(\vec{x} + \delta \vec{v}) - f(\vec{x})}{\delta} = \sum_i v_i \frac{\partial f}{\partial x_i}$

Can show $D_{\vec{v}}\left(\frac{\partial E}{\partial \theta_i}\right) = [H\vec{v}]_i$

- In AutoDiff:

1. Forward pass, ordering nodes

2. Backward pass, calculate t_n

3. Directional derivative forward pass

$$D_{\vec{v}}(f_i) = \sum_k \frac{\partial f_i}{\partial \pi_k} D_{\vec{v}}(\pi_k) \quad (\vec{\pi}_k \text{ are parameters})$$

4. Directional derivative reverse pass:

$$D_{\vec{v}}(t_n) = \sum_{c \in h(n)} \left(\frac{\partial \xi_c}{\partial g_n} D_{\vec{v}}(t_c) + t_c \frac{\partial^2 \xi_c}{\partial g_n^2} D_{\vec{v}}(f_n) \right)$$

7c: Neural Net Initialisation

- First, scale inputs to 0 mean, unit variance

$$\vec{x}_i^n \leftarrow \frac{\vec{x}_i^n - \mu_i}{\sigma_i}$$

- Aim to keep output 0 mean, unit variance \leftarrow to keep error $O(1)$

\Rightarrow Sample weights from 0 mean, $1/\text{size of input } \vec{x}$ variance

can do for multiple layers

- Decorrelating inputs

approximate Hessian, & apply

$$\vec{x}^n \leftarrow S^{-1} U^T \vec{x}^n, \text{ where } X X^T = U S^2 U^T$$

(ie SVD of X)

make approximate Hessian diagonal & can be useful to some optimisation

8a: Recurrent Networks

- Gradient decay/exlosion:

- for a very deep network, gradients grow exponentially

- w/ no. layers \Rightarrow can explode or decay to 0

- need special architectures to deal w/ this

- Memory decay:

- information typically lost over time steps in traditional RNN

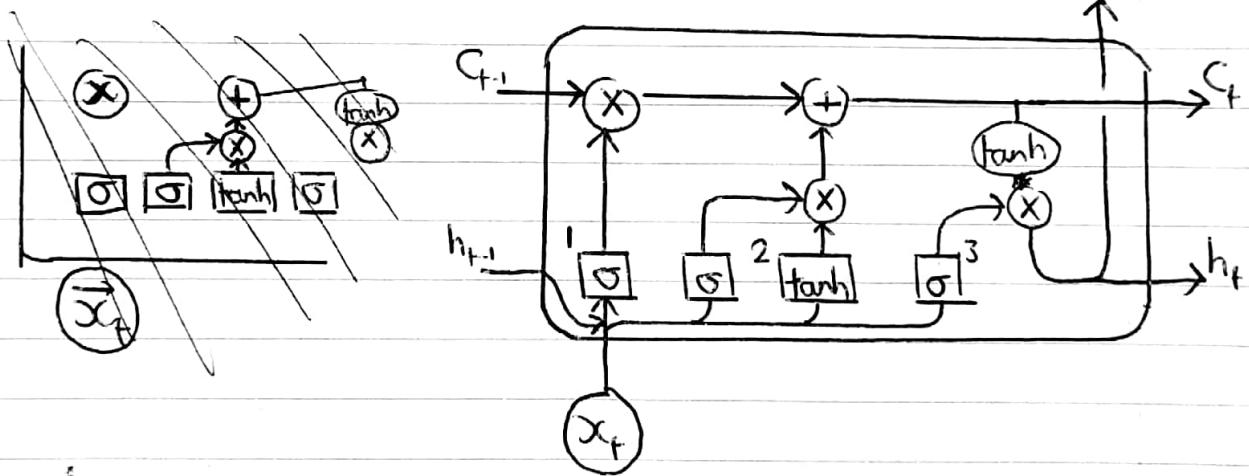
- need some form of 'running memory', $m(t)$

- can be 'gated' using $\alpha(t) = \sigma(W_{in}x(t) + B_{in})$

- $m(t) = \alpha(t)x(t) + (1-\alpha(t))m(t-1)$

- works well when α close to 1; gradient decays slowly

- LSTM



- top is memory

- 1: forget gate; allow for forgetting memory

- 2: input gate; update memory

- 3: output gate; decide what bit of memory to output

- GRU:

- similar, but h_t acts as memory & is fully exposed as output

8b: Unsupervised Dimension Reduction

- In high dimensional data with some structure, it will tend to lie close to some lower dimensional manifold

- Principle Component Analysis:

• Full rank representation: $\vec{x} = \sum_{d=1}^D y_d \vec{b}^d$, $y_d \in \mathbb{R}^0$

• Reduced rank approximation: $\vec{x} \approx \sum_{d=1}^M (y_d \vec{b}^d) + \vec{c}$, i.e. $\vec{x} = \vec{c} + B\vec{y}$

note: to give unique solution,
ensure all \vec{b}^d orthogonal & unit length

• aim to minimise $E(B, Y, \vec{c}) = \sum_n \sum_d (x_d^n - \tilde{x}_d^n)^2$

• can show $\vec{c} = \frac{1}{N} \sum_n \vec{x}^n$, and \vec{b}^d are eigenvectors of sample covariance matrix, $\frac{1}{N-1} \sum_n (\vec{x}^n - \vec{c})^T (\vec{x}^n - \vec{c})$

• can reduce dimensionality by only taking top M eigenvectors (by eigenvalue)

• Eigen-spectrum:

• ~~total~~ squared error of approximation

$$E = (N-1) \sum_{n=M+1}^N \lambda_n$$

i.e., sum neglected eigenvalues

• ~~Rotational invariance~~

• PCA Maths:

$$\text{minimize } E(B, Y) = \sum_{n,d} (x_d^n - \sum_{j=1}^M y_j^n b_j^d)^2 = \text{Trace}((X - BY)^T (X - BY))$$

$$\frac{\partial E}{\partial y_k} = 0 \Rightarrow y_k^n = \sum_d b_d^n x_d^n$$

$$X = (\vec{x}^1, \dots, \vec{x}^N)$$

$$B^T B = I \quad (\text{i.e. basis vectors are orthogonal \& unit length})$$

$$Y = B^T X$$

$$\Rightarrow E(B) = \text{trace}(X^T (I - B B^T)^2 X) = \text{trace}(X X^T (I - B B^T))$$

Sample Covariance

$$= [\text{trace}(S) - \text{trace}(S B B^T)] (N-1)$$

To ensure $B^T B = I$, use Lagrange multipliers L

$$\Rightarrow \text{minimize } -\text{trace}(S B B^T) + \text{trace}(L(B^T B - I))$$

~~at optimum~~

\Rightarrow at optimum $S\beta = BL$

• enforce L to be diagonal \Rightarrow cols of B are corresponding eigenvectors of S

• Now, $\text{trace}(SBB^T) = \text{trace}(L)$

$$\Rightarrow \frac{1}{N-1} E(B) = E(B) = \text{Tr}(S) - \text{Tr}(L) = \sum_{d=1}^D \lambda_d - \sum_{d=1}^M \lambda_d = \sum_{d=M+1}^D \lambda_d$$

Squared residual error!

• PCA w/ missing data:

• some x_i^n missing; aim to infer them

• first, find B & Y by only minimising error wrt only known x_i^n

• need to iterate B for fixed Y & vice versa, decreasing loss until minimum reached

• can be used e.g. for Movie recommendation from a matrix of liked/disliked movies

• Matrix decompositions

$$X \approx B Y$$

$(D \times N) \quad (D \times M) \quad (M \times N)$

Data basis weights or components

(PCA is a type of)
(this, where $M \leq D$)

• Under-complete: $M < D$, e.g. PCA

• Y is lower dimensional approximate representation of X

• Over-complete: $M > D$

• need to introduce additional constraints

• can lead to sparse representations of data

8b: Unsupervised Dimension Reduction

• Non-Negative Matrix Factorisation (NMF)

$$X_{ij} \approx \sum_k B_{ik} Y_{kj}, \text{ where } X_{ij}, B_{ik} \text{ & } Y_{kj} \text{ all +ve}$$

• e.g. X columns are images, B are 'basis' vectors for these images & Y weights for basis vectors

• Training:

• use an EM algorithm

$$\text{• aim to approximate } P(x, y) \approx \underbrace{\sum_z P(x|z) P(y|z) P(z)}_{\hat{P}(x, y)}$$

• M: fix $P(x|z)$, $P(y|z)$, $P(z)$

• find $P(z)$ to minimise $KL(P(x, y), \hat{P}(x, y))$

• E: fix $P(z)$

• similarly, update $P(x|z)$ & $P(y|z)$

• E: update $q(z|x, y)$ used in M step

• guaranteed to increase likelihood \Rightarrow converges

+ interpretable basis vectors

9a: Independent Component Analysis

- Find coordinate system where coords are independent
 - compared w/ PCA, where enforced to be orthogonal
 - can be better than PCA, for eg, finding directions along which data was generated (which not necessarily orthogonal)
- Train:
 - maximize likelihood
$$L(A) = L(B) = N \log \det(B) + \sum_{n,i} \log P(B \vec{v}^n)_i$$

where $\vec{v}^1 \dots \vec{v}^N$ are points & A being our ICA estimate

 - derive a gradient ascent rule (H₀ can apply Newton's method to get natural gradient update)

9b: Supervised Dimension Reduction

- Aim to find linear projection of data to lower dimensional space that allows us best classification after performance
- Fisher's Linear Discriminant:
 - model data as produced by Gaussian for each class
 - projected distribution still Gaussian
 - project such that there is minimal overlap

$$p(\vec{x}_{oi}) = N(\vec{x}_i | \vec{m}_i, S_i)$$

\vec{m}_i = ^{Sample} Mean of group 1
 S_i = ^{Sample} covariance

$$y_i = \vec{w}^T \vec{x}_i$$

Find \vec{w} such that minimal overlap in y 's

$$\min \left(\frac{(\mu_1 - \mu_2)^2}{\pi_1 \sigma_1^2 + \pi_2 \sigma_2^2} \right) = \frac{\vec{w}^T (\mu_1 - \mu_2) (\mu_1 - \mu_2)^T \vec{w}}{\vec{w}^T (\pi_1 S_1 + \pi_2 S_2) \vec{w}} = \frac{\vec{w}^T A \vec{w}}{\vec{w}^T B \vec{w}}$$

$$\Rightarrow \vec{w} \propto B^{-1} (\vec{m}_1 - \vec{m}_2)$$

• If $N_1 + N_2 < D$, B not invertible (e.g. features which never change across all \vec{x} can also make objective ill defined)

- Canonical Variates:

• Generalize Fisher's method to > 1 dimension & > 2 classes

$$p(\vec{x}) = N(\vec{x} | \vec{m}_c, S_c),$$

$$\vec{y} = \vec{w}^T \vec{x} \Rightarrow p(\vec{y}) = N(\vec{y} | \vec{w}^T \vec{m}_c, \vec{w}^T S_c \vec{w})$$

• Compute $\underline{A} = \sum_c N_c (\vec{m}_c - \vec{m})(\vec{m}_c - \vec{m})^T$

class count \downarrow class c mean \downarrow mean of all data \downarrow

$$\underline{B} = \sum_c N_c \underline{S}_c$$

• Let $\tilde{\underline{B}}^T \tilde{\underline{B}} = \underline{B}$ (assuming \underline{B} invertible) (\underline{W} orthonormal)

• Maximise $F(\underline{W}) = \frac{\text{tr}(\underline{W}^T \tilde{\underline{B}}^{-T} \underline{A} \tilde{\underline{B}}^{-1} \underline{W})}{\text{tr}(\underline{W}^T \underline{W})} = \frac{\text{tr}(\underline{W}^T \tilde{\underline{B}}^{-T} \underline{A} \tilde{\underline{B}}^{-1} \underline{W})}{\underline{W}^T \underline{W}}$

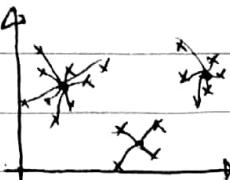
• $\underline{C} = \underline{E} \underline{\Lambda} \underline{E}^T$ (PSD & symmetric)

$$F(\underline{W}) = \text{tr}(\underline{W}^T \underline{E} \underline{\Lambda} \underline{E}^T \underline{W})$$

[eigenvalues of $\underline{B} \underline{C}$]

\Rightarrow maximised by $\underline{W} = [\vec{e}_1, \dots, \vec{e}_L]$, where

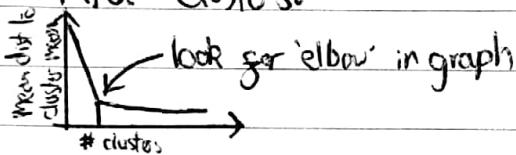
9c: Clustering



- **K-Means:**

- minimise sum of squared distances of points from cluster mean
- Initialise means randomly
- until converge:
 - assign points to closest cluster centre
 - recalculate cluster centres (mean of points in cluster)
- decreases loss each iteration \Rightarrow will converge
- non-deterministic; repeat with different initialisations

- **No. clusters:**



- **Limitations:**

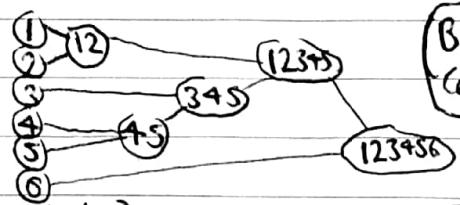
- only 1 cluster per each datapoint
- sensitive to outliers (K-medians better)
- requires spherical clusters of equal density
- scales poorly with dimensionality
- sensitive to dimension scale (~~all~~* should be rescaled)
- requires numerical data; encode categorical w/ one-hot vectors
- how to deal with missing geofields?

* each dimension

- **Hierarchical K-means:**

- can do bottom up
(merge from single points)
- or top down (remove

'outsider' from least cohesive cluster)



- **Spectral Clustering:**

- model data as Markov chain, with

$$p(x'|x) = \frac{A_{x,x'}}{\sum_{x''} A_{x,x''}}, \text{ where } A_{i,j} \text{ is similarity of } i \text{ & } j$$

- find equilibrium distribution, & can cluster this with K-means

result is very dependent on choice of this!
Bad?

• Gaussian Mixture Models:

- model data generation; first sample cluster, then sample point from that cluster

$$P(v) = \sum_h P(v|h)P(h)$$

cluster h ('hidden')

point v ('visible' or 'observable')

- Learn parameters by maximising likelihood (use EM algorithm)
- Cluster by computing $\hat{p} P(h|v) \propto P(v|h)P(h)$ ($\hat{p} = \max P(h|v)$)
- $P(v|h)$ need not be Gaussian; could e.g. be Bernoulli (for clustering binary (e.g. some features from categorical distribution))

$$L = \sum_n \log \sum_c P(c) \times \frac{1}{\sqrt{det(2\pi \underline{\underline{S}}_c)}} \exp\left(-\frac{1}{2} (\underline{x}^n - \underline{m}_c)^T \underline{\underline{S}}_c^{-1} (\underline{x}^n - \underline{m}_c)\right)$$

- $\underline{\underline{S}}_c$ must be symmetric, PSD
- if $\underline{m}_c = \underline{x}^n$ for some c, n , as $\underline{\underline{S}}_c \rightarrow 0$, $P(\underline{x}^n | \underline{m}_c, \underline{\underline{S}}_c)$ becomes infinite \Rightarrow need to limit how small $\underline{\underline{S}}_c$ ('width' of Gaussian) can be

• k-means:

- if constrain all $\underline{\underline{S}}_c$ to be $\sigma^2 \underline{\underline{I}}$, as $\sigma^2 \rightarrow 0$, GMM becomes equivalent to k-means
- $\sigma^2 > 0$ gives 'softer' clustering than k-means
- Gaussian not too resistant to outliers; use t-distribution

-  Find number of pure clusters by cross-validation (likelihood of test data)

- Easy to train using EM

10: Visualisation

- Aim to find low dimensional representation (2 or 3) of data
• no correct or perfect method; all lose information

Methods:

- PCA & Sammon 'mappings'; old, less preferred now
- Autoencoders; good, very expensive to train

Stochastic Neighbour Encoding:

- aim to transform \mathbb{X} to \mathbb{Y} (high \rightarrow low dimensional)

$$p_{j|i} = \frac{\exp(-(\mathbb{x}_i - \mathbb{x}_j)^2 / (2\sigma^2))}{\sum_{j \neq i} \exp(-(\mathbb{x}_i - \mathbb{x}_j)^2 / (2\sigma^2))} \quad (p_{i|i} = 0)$$

$$q_{i|i} = \frac{\exp(-(\mathbb{y}_i - \mathbb{y}_i)^2)}{\sum_{j \neq i} \exp(-(\mathbb{y}_i - \mathbb{y}_j)^2)} \quad (q_{i|i} = 0)$$

do so
by gradient
descent;
treat \mathbb{y}_i 's
as parameters
to optimise

- aim to find \mathbb{Y} such that q similar to p
- minimise $\sum_i \text{KL}(p_i \| q_i) = \sum_i p_{i|i} \log \frac{p_{i|i}}{q_{i|i}}$

Problems:

- KL divergence not symmetric; leads to good local structure but poor global

- Gaussian $p_{i|i}$ = distant points have too little effect

t-SNE:

- use 'symmetric' loss

$$\text{KL}(p \| q) = \sum_{i,j} p_{i,j} \log \frac{p_{i,j}}{q_{i,j}}$$

$$\left(\text{where } p_{i,j} = \frac{p_{i|i} + p_{j|j}}{2N}, p_{i|i} = 0 \right)$$

- let $q_{i,j}$ be a Student t-distribution

$$q_{i,j} = \frac{(1 + (\mathbb{y}_i - \mathbb{y}_j)^2)^{-1}}{\sum_{i \neq j} (1 + (\mathbb{y}_i - \mathbb{y}_j)^2)^{-1}}, q_{i|i} = 0$$

\Rightarrow very distant points similar contribution to somewhat
distant points

- Scaling:

- $O(N^2)$ to compute t-SNE objective; so $O(N^2)$ per iteration of gradient descent!
- Alternative:
 - select a few landmark points, i' , randomly
 - Compute ~~PS~~ $p_{j|i'}$ for all pairs of landmarks by doing many markov chains from i' until it visits another landmark, j
 - now do t-SNE on smaller $p_{j|i'}$
 - expensive to compute $p_{j|i'}$, but only needs to be done once

11: Fast Nearest Neighbours Calculation

- Finding nearest neighbour naively is $O(DN)$

- Orchard:

- good when distance computation expensive (eg high D)

- Triangle inequality:

$$|\underline{x} - \underline{y}| \leq |\underline{x} - \underline{z}| + |\underline{z} - \underline{y}| \quad (\text{ie faster to go to } \underline{y} \text{ directly})$$

- Metric:

- $d(\underline{x}, \underline{y}) \geq 0, d(\underline{x}, \underline{y}) = 0 \Rightarrow \underline{x} = \underline{y}$,

- $d(\underline{x}, \underline{y}) = d(\underline{y}, \underline{x}) \quad (\text{& } d(\underline{x}, \underline{y}) \leq d(\underline{x}, \underline{z}) + d(\underline{z}, \underline{y}))$

- Euclidean distance is a metric

- Exploiting triangle inequality:

- if $d(\underline{x}, \underline{y}) \leq \frac{1}{2}d(\underline{x}, \underline{z}) \Rightarrow d(\underline{x}, \underline{y}) \leq d(\underline{x}, \underline{z})$
 \Rightarrow don't need to calculate $d(\underline{x}, \underline{z})$

- Algorithm:

$O(DN^2)$
time!
 $O(N^2)$
space!

- precompute all $d_{i,j} = \text{dist}(\underline{x}_i, \underline{x}_j)$

- for each i , compute ordered list of \underline{x}_j closest to \underline{x}_i

- for query point q :

- pick \underline{x}_i , initial guess at q 's NN

- go along \underline{x}_i 's list, finding closest point to q

- ~~can stop early~~

- if $d(\underline{x}_j, q) < d(\underline{x}_i, q)$, go to \underline{x}_j 's list

- if $d(\underline{x}_j, q) \geq \frac{1}{2}d(\underline{x}_i, q)$ can stop early!

- Approximating & Eliminating Search Algorithm (AES):

- again, precompute distance matrix

- $d(q, \underline{x}^j) \geq d(q, \underline{x}^i) - d(\underline{x}^i, \underline{x}^j) \quad (\text{triangle inequality})$

- \Rightarrow compute set $J = \{\text{points for which } d(q, \underline{x}^i) \text{ calculated}\}$

- $d(q, \underline{x}^j) \geq \max_{i \in J} \{d(q, \underline{x}^i) - d(\underline{x}^i, \underline{x}^j)\}$

$O(M(N-M))$

\uparrow
repeat

$M \leq N$ times

- eliminate all j where this bound is greater than

$$\min d(q, \underline{x}^i)$$

$O(N^2)$ space

- **Buoys:**

- Storing full distance matrix can be prohibitive
- instead, store distances to set of buoys, $\underline{b}_1, \dots, \underline{b}_m$
- Pre-elimination:

- $d(q, \underline{x}^n)$

$$\max_m \{d(q, \underline{b}^m) - d(\underline{b}^m, \underline{x}^n)\} \leq d(q, \underline{x}^n) \leq \min_m \{d(q, \underline{b}^m) + d(\underline{b}^m, \underline{x}^n)\}$$

- eliminate any points where lower bound is larger than smallest upper bound

- then feed remainder to AESA or Orchard

AKA
Linear
AES

- AESA with buoys:

- order by lower bounds

- evaluate, if $d(q, \underline{x}^n) \leq L_{\text{max}}$, terminate; \underline{x}^n is NN

- Orchard w/ buoys:

- worse than Linear-AESA

- k-d tree:

- binary tree, recursively dividing space after through median, cycling through dimensions to divide along

- Finding nearest neighbour:

- traverse tree with q to root or leaf, this is best guess

- traverse up tree:

- check parent, if better ~~neighbor~~ then set this as best guess

- check if need to traverse down other to other child of parent (hypersphere of radius = current best guess intersecting w/ hyperplane of parent = need to check other child)

