

Decide for  $C_i$ , if

$$p(C_i|x) > p(C_{i+1}|x) \Rightarrow \frac{p(x|C_i)}{p(x|C_{i+1})} > \frac{p(C_{i+1})}{p(C_i)} \quad \text{decision threshold } \theta$$

Likelihood-ratio test for  $k$  class

$$\frac{p(x|C_k)}{p(x|C_j)} > \frac{p(C_j)}{p(C_k)} \quad \forall j \neq k$$

Minimize the expected loss:

$$E(L) = \sum_k \sum_j \int_{R_j} L_{kj} p(x, C_k) dx$$

Generalisation with loss function

We can formalize this by introducing a loss matrix  $L_{kj}$

$$L_{\text{cancer}} = \begin{pmatrix} \text{(Truth)} & \text{cancer} \\ \text{normal} & \text{normal} \end{pmatrix} \begin{pmatrix} 0 & 1000 \\ 1 & 0 \end{pmatrix}$$

$2$  class:  $C_1, C_2$ ,  $2$  decision:  $\alpha_1, \alpha_2$

Loss function:  $L(\alpha_j | C_k) = L_{kj}$

Expected loss ( $=$  risk  $R$ ) for two decisions:

$$E_{\alpha_1}(L) = R(\alpha_1|x) = L_{11} p(C_1|x) + L_{12} p(C_2|x)$$

$$E_{\alpha_2}(L) = R(\alpha_2|x) = L_{21} p(C_1|x) + L_{22} p(C_2|x)$$

Discriminant function

$$y_1(x), \dots, y_K(x) \leftarrow \begin{cases} p(x|C_i) & \text{if } y_K(x) > y_j(x) \\ p(x|C_j) & \text{if } y_j(x) > y_i(x) \quad \forall i \neq j \end{cases}$$

$$y_K(x) = \log p(x|C_K) + \log p(C_K)$$

Goal: decide such that expected loss is minimized  
i.e. decide  $\alpha_i$ , if  $R(\alpha_i|x) > R(\alpha_j|x)$

$$N(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$$

$$N(x|\mu, \Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right\}$$

Data:  $x_1, \dots, x_n$  estimate:  $p(x)$

- Method:
- Parametric representation
  - Non parametric representation
  - Mixture models

Given data  $X = \{x_1, x_2, \dots, x_N\}$ . Parametric form of the distribution with parameters  $\theta$ .

E.g. for Gaussian distribution  $\theta = (\mu, \Sigma)$

Learning: Estimate parameter  $\theta$

Likelihood of  $\theta$

Probability that the data  $x$  have indeed been generated from a probability density with parameter  $\theta$ .  $L(\theta) = p(x|\theta)$

MLE - Computation of the likelihood

• Single data point:  $p(x_1|\theta)$

• Assumption: all data points are independent

$$L(\theta) = p(x|\theta) = \prod_{i=1}^N p(x_i|\theta)$$

$$\text{Log likelihood } \ell(\theta) = -\ln L(\theta) = -\sum_{n=1}^N \ln p(x_n|\theta)$$

Estimate of the parameters  $\theta$

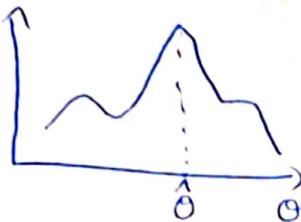
• Maximize the likelihood

• Minimize the negative log-likelihood

likelihood:  $L(\theta) = p(x|\theta) = \prod_{i=1}^N p(x_i|\theta)$

• we want to obtain  $\hat{\theta}$  such that  $L(\hat{\theta})$  is maximized.

$p(x|\theta)$



Minimizing the log-likelihood.

$$\frac{\partial E(\theta)}{\partial \theta} = -\frac{1}{\theta} \sum_{i=1}^N \ln p(x_i|\theta) = -\sum_{n=1}^N \frac{\frac{\partial}{\partial \theta} p(x_n|\theta)}{p(x_n|\theta)} = 0$$

Log likelihood for normal distribution.

$$E(\theta) = -\sum_{i=1}^N \ln p(x_i|\mu, \sigma^2) = -\sum_{n=1}^N \ln \left( \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(x_n - \mu)^2}{2\sigma^2} \right\} \right)$$

$$\text{Thus } \hat{\mu} = \frac{1}{N} \sum_{n=1}^N x_n \quad \hat{\sigma}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \hat{\mu})^2 \quad \hat{\theta} = (\hat{\mu}, \hat{\sigma}^2)$$

$$E(\mu_{ML}) = \mu \quad E(\sigma_{ML}^2) = \left( \frac{N-1}{N} \right) \sigma^2$$

MLE

→ underestimate the variance of the distribution.

→ overfits the observed data.

Probability that  $x$  falls into small region  $R$ .

$$P = \int_R p(y) dy \quad P \approx p(x)V \quad P = \frac{K}{N} \Rightarrow p(x) \approx \frac{K}{NV}$$

Kernel method      K nearest neighbour.

Parzen window

$$k(u) = \begin{cases} 1 & |x_i| \leq \frac{1}{2}h, i=1..D \\ 0 & \text{else} \end{cases}$$

$$K = \sum_{n=1}^N k(x-x_n) \quad V = \int k(u) du = h^D$$

$$\text{Probability density estimate} \quad p(x) \approx \frac{K}{NV} = \frac{1}{Nh^D} \cdot \sum_{n=1}^N k(x-x_n)$$

Gaussian kernel

$$k(u) = \frac{1}{(2\pi h^2)^{D/2}} \exp\left\{-\frac{u^2}{2h^2}\right\}$$

$$K = \sum_{n=1}^N k(x-x_n) \quad V = \int k(u) du = 1$$
$$p(x) \approx \frac{1}{N} \sum_{n=1}^N \frac{1}{(2\pi)^{D/2} h} \exp\left\{-\frac{\|x-x_n\|^2}{2h^2}\right\}$$

K-nearest neighbour.

$$p(c_j|x) = \frac{p(x|c_j)p(c_j)}{p(x)} \quad p(x) \propto \frac{1}{NV} \quad p(x|c_j) \approx \frac{k_j}{NJV} \rightarrow$$
$$p(c_j|x) \approx \frac{k_j}{NJV} \cdot \frac{N_j}{N} \cdot \frac{NV}{1C} = \frac{k_j}{1C}$$

$$p(x|x) = \int \frac{p(x|\theta)L(\theta)p(\theta)}{\int p(x|\theta)p(\theta)d\theta} d\theta \quad \boxed{p(c_j) \approx \frac{N_j}{N}}$$

## Mixture models

- hard clustering
  - element either belong to cluster or not
- soft clustering
  - soft assignment

- Mixture models are probabilistically doing soft clustering.

- each cluster: a generative model (Gaussian or multinomial)
- parameters e.g. mean / covariance are unknown

⇒ EM algo used to find the parameters for K "sources"

-  $K=2$ , Gaussian with unknown  $\mu, \sigma^2$

- If I know the source  $i$  can find  $\mu_i, \sigma^2$

If we know the source of the gaussian ( $\mu, \sigma^2$ ) we can guess more likely from "a" or "b" using bayes

$$p(b|x_i) = \frac{p(x_i|b) p(b)}{\int_{\Omega} p(x_i) \exp\left(-\frac{(x_i - \mu_b)^2}{2\sigma^2}\right) d\mu}$$

we can use EM

- start with random Gaussian  $(\mu_a, \sigma_a^2), (\mu_b, \sigma_b^2)$

- for each point:  $p(b|x_i) = \text{does it look like it came from } b?$

- adjust  $(\mu_a, \sigma_a^2)$  and  $(\mu_b, \sigma_b^2)$  to fit points assigned to them

↳ E-step: soft assignment

$$\gamma_j(x_n) \leftarrow \frac{\pi_j N(x_n | \mu_j, \Sigma_j)}{\sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k)}$$

↳ M-step: re-estimate the parameters

$$\mu_j \leftarrow \frac{1}{N_j} \sum_{n=1}^N \gamma_j(x_n) x_n$$

$$\Sigma_j \leftarrow \frac{1}{N_j} \sum_{n=1}^N \gamma_j(x_n) (x_n - \mu_j)^T (x_n - \mu_j)$$

## GMM or MOG

$$p(x|\theta) = \sum_{j=1}^K p(x|\theta_j) p(j)$$

$$p(x|\theta_j) = N(x | \mu_j, \Sigma_j) = \frac{1}{\sqrt{2\pi\Sigma_j}} \exp\left(-\frac{(x - \mu_j)^2}{2\Sigma_j}\right)$$

$$p(j) = \pi_j \quad 0 \leq \pi_j \leq 1$$

prior of component "j"

$$\text{MLE} \rightarrow \hat{\theta} = \ln L(\theta) = \sum_{n=1}^N \ln p(x_n|\theta)$$

$$\mu_j = \frac{\sum_{n=1}^N \gamma_j(x_n) x_n}{\sum_{n=1}^N \gamma_j(x_n)}$$

$$\begin{aligned} \gamma_j(x_n) &= \frac{\pi_j N(x_n | \mu_j, \Sigma_j)}{\sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k)} \\ &\text{Responsibility of component } j \text{ for } x_n \end{aligned}$$

## K-mean

if optimise the following

$$J = \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \|x_n - \mu_k\|^2$$

$$\gamma_{nk} = \begin{cases} 1 & \text{if } k = \arg\min_{l \in \{1, \dots, K\}} \|x_n - \mu_l\|^2 \\ 0 & \text{otherwise} \end{cases}$$

1. Initiation: pick K arbitrary centroids (cluster means)
2. Assign each sample to the closest centroid.
3. Adjust the centroid to be the means of the samples assigned to them
4. Go to step 2 (until no change)

## Pros

- simple & fast to compute
- converges to local minimum of within-cluster squared error

## Problem

- settings K?
- sensitive to initial centres
- sensitive to outliers
- Detect only spherical clusters

## GMM

- very general, can represent any distribution
- once trained, very fast evaluation
- can be updated online

## Problem

- some numerical issues in implementation
- need to apply regularisation to avoid singularities
- EM of MOG is computationally expensive
  - especially for high-dimensional problems
  - slower convergence than K-means
  - result sensitive to initialisation
- Need to select number of mixture components K.

• Discriminant functions  
 $y_1(x), \dots, y_K(x)$

$$y_k(x) > y_j(x) \quad \forall j \neq k$$

$$y_k(x) = \log p(x|C_k) + \log P(C_k)$$

$$y_1(x) > y_0(x)$$

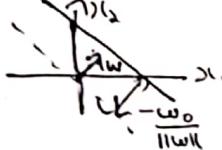
$$y(x) > 0$$

### 2-class problem

$y(x) > 0$ : decide for  $C_1, C_0$

$$y(x) = w^T x + w_0 = \sum_{i=0}^n w_i x_i$$

$$y(x) = 0 \rightarrow \text{decision boundary } w_0 = 0$$



### One-vs-all classifiers

heuristic method for using binary classification algorithm for multi-class.

Ex: R, B, G

Binary classification:  $R \vee [B|G]$

" " 2:  $R \vee [R|G]$

" " 3:  $G \vee [R,B]$

It requires one model to be created for each class.

### one vs one multiclass

Ex: R, B, G, Y

Binary class 1: R vs B

2: R vs G

3: R vs Y

4: B vs G

5: B vs Y

6: G vs Y

No. of binary datasets  $\frac{N \times (N-1)}{2}$

### fc class

$$y_{1k}(x) = w_{1k}^T x + w_{1k0}, \quad k=1, \dots, K$$

$$y_{jk}(x) = w_{jk}^T x$$

$$\tilde{w} = [\tilde{w}_{11} \dots \tilde{w}_K]$$

$$\text{target } t = [t_1, \dots, t_K]^T$$

$$Y(\tilde{x}) = \tilde{x} \tilde{w}$$

$$\tilde{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} \quad T = \begin{bmatrix} t_1^T \\ \vdots \\ t_K^T \end{bmatrix}$$

$$\tilde{x} \tilde{w} - T$$

### Least-squares

Sum of square error

$$E(w) = \frac{1}{2} \sum_{m=1}^M \sum_{k=1}^K (w_{1k}^T x_m - t_{1k})^2$$

### Multiclass

$$E(w) = \frac{1}{2} \operatorname{Tr} \{ (\tilde{x} \tilde{w} - T)^T (\tilde{x} \tilde{w} - T) \}$$

$$\frac{\partial E(w)}{\partial w} = \tilde{x}^T (\tilde{x} \tilde{w} - T) \Rightarrow y(x) = \tilde{w}^T \tilde{x} = T^T (\tilde{x}^T \tilde{w}) \tilde{x}$$

pseudo inverse

~~Problem~~

~~Solve~~

Take K linear functions of form  $y_{1k}(x) = w_{1k}^T x + w_{1k0}$  decision boundary  $C_k$

iff  $y_k > y_j, \forall j \neq k$

~~Problem~~ 1-of-K coding scheme

$$t_m = (0, 1, 0, 0, 0)^T$$

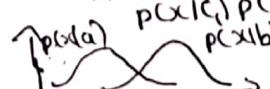
$$(w_{1j} - w_{1k})^T x + (w_{1k0} - w_{1j0}) = 0$$

### Least square

Very sensitive to outliers  
 Least square corresponds to Maximum likelihood under the assumption of a Gaussian conditional distribution.  
 activation function may be nonlinear.

$$y(x) = g(C w^T x + w_0)$$

$$p(C, t|x) = \frac{1}{1 + \frac{p(x|C_a)p(t|C_a)}{p(x|C_b)p(t|C_b)}} = \frac{1}{1 + \exp(-a)} = g(a) \rightarrow \text{logistic sigmoid activation function}$$



Scanned by CamScanner

## LDA II

If data is not linearly separable, a linear decision boundary is still optimal.  
Choice of discriminant function is based on

- Prior knowledge - empirical comparison of alternative method.

### Generalized linear discriminants

Transform vector  $x$  with  $M$  nonlinear basis function  $\phi_j(x)$ : This basic function

$$y_k(x) = \sum_{j=1}^M w_{kj} \phi_j(x) + w_{k0} = y_k(x; w) \quad \text{allows nonlinear boundaries.}$$

Single-layer network:  $\phi_j$  is fixed, only weight  $w$  is learned.

Multi-layer network: both  $w$  and  $\phi_j$  are learned.

GD (Gradient descent)  $\rightarrow$  feature  
sample:  $sc = [x_1, x_2, \dots, x_n] \rightarrow$  single data point  $y$  (MLE)  
 $x \in \mathbb{R}^d \quad y \in \{0, 1\}$  [Training: find  $\theta$  that maximize seeing training data]

$D = \{(x_i, y_i)\}_{i=1}^n$ : we want to learn  $\theta = [\theta_0, \theta_1, \dots, \theta_M]$

$$\hat{\theta}_{MLE} = \arg \max_{\theta} P(D; \theta) = \arg \max_{\theta} P(Y|X; \theta) = \arg \max_{\theta} \prod_{i=1}^n P(y_i|x_i; \theta)$$

For logistic  $P(y_i|x_i; \theta) = \text{Function}(\theta^T x_i) \rightarrow$  sigmoid function  $\sigma(\theta^T x_i) \rightarrow$  link function

$$P(y_i=1|x_i; \theta) = \sigma(\theta^T x_i) \rightarrow L(\theta) = \max_{\theta} \prod_{i=1}^n P(y_i|x_i; \theta) = \max_{\theta} \prod_{i:y_i=1} \sigma(\theta^T x_i) \prod_{i:y_i=0} (1 - \sigma(\theta^T x_i))$$

$$L(\theta) = \prod_{i=1}^n \sigma(\theta^T x_i)^{y_i} \times (1 - \sigma(\theta^T x_i))^{1-y_i} \quad \text{(negative log likelihood)}$$

$$l(\theta) = \max_{\theta} \sum_{i=1}^n y_i \log(\sigma(\theta^T x_i)) + (1-y_i) \log(1 - \sigma(\theta^T x_i)) \quad \text{NLL}(\theta) = -l(\theta)$$

$$NLL(\theta) = \min_{\theta} \sum_{i=1}^n y_i \log(\sigma(\theta^T x_i)) + (1-y_i) \log(1 - \sigma(\theta^T x_i)) \quad \text{① we need to minimize the log likelihood} \quad \text{loss function}$$

NLL( $\theta$ ) = loss function. Here it is GD.  
Optimizer used to find  $\theta$  that minimizes loss function. Here it is based on Gaussian

1st order optimizers: GD, require first order derivative, Jacobian matrix, which is based on Newton-Raphson, 2nd order derivative, Hessian matrix, which is based on Hessian is more.

GD  $\theta = \theta + \alpha \frac{\delta NLL}{\delta \theta^T} \rightarrow$  loss change / unit change in weights.  $\frac{d(\sigma(x))}{dx} = \sigma(x)(1 - \sigma(x))$

$$\theta = \theta + \alpha \frac{\delta NLL}{\delta \theta^T} = \sum_{i=1}^n (y_i - \sigma(\theta^T x_i)) \cdot x_i$$

$$\text{chain rule to ①} \rightarrow \frac{\delta NLL}{\delta \theta^T} = \sum_{i=1}^n (y_i - \sigma(\theta^T x_i)) x_i, \hat{\theta}_{MLE} = \theta$$

GD  $\rightarrow$  iterate over  $m$  iterations  $\theta = \theta + \alpha \sum_{i=1}^m (y_i - \sigma(\theta^T x_i)) x_i$

$$P(y=1|x) = \sigma(\theta^T x; \theta = \hat{\theta}_{MLE}) = \frac{1}{1 + e^{-\hat{\theta}_{MLE}^T x}}$$

$$\text{error function} E(w) = \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K (y_{nk}(x_n; w) - t_{kn})^2 = \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K \left( \sum_{j=1}^M w_{kj} \phi_j(x_n; w) - t_{kn} \right)^2$$

$$w_{kj}^{(t+1)} = w_{kj}^{(t)} - \eta \frac{\partial E(w)}{\partial w_{kj}} \Big|_{w^{(t)}} \quad \frac{\partial E(w)}{\partial w_{kj}} = (y_{kj}(x_n; w) - t_{kn}) \phi_j(x_n)$$

$$w_{kj}^{(t+1)} = w_{kj}^{(t)} - \eta \delta_{kn} \phi_j(x_n) \quad \delta_{kn} = y_{kj}(x_n; w) - t_{kn}$$

with activation function  $y_{kj}(x_n; w) = \sigma(\theta^T x_n; w)$   $\{ \delta_{kn} = \frac{\partial g(\theta^T x_n)}{\partial w_{kj}} (y_{kj}(x_n; w) - t_{kn}) \}$

$$\frac{\partial E(w)}{\partial w_{kj}} = ( ) \Rightarrow w_{kj}^{(t+1)} = w_{kj}^{(t)} - \eta \delta_{kn} \phi_j(x_n)$$

Probabilistic linear discrimination and gaussian method of modelling.

No. of parameters for gaussian  $\rightarrow$  mean:  $M$  + covariance  $M(M+1)$   $\xrightarrow{\text{dimension}}$  total  $M(M+1)/2 + 1$  parameters

logistic regression parameters: just the values of  $w \Rightarrow M$  parameters

Consider dataset,  $\{(\phi_n, t_n)\}$  with  $n=1, \dots, N$ , where  $\phi_n = \phi(x_n)$  and  $t_n \in \{0, 1\} t = \{t_1, \dots, t_N\}$

### Logistic regression

Classify the weather as rainy, sunny and partially cloudy.

$$P(y=1|x) = p(x) \quad x \in \mathbb{R} \quad p(x) \in [0, 1], \quad p(x) = \frac{1}{1+e^{-\beta x}}$$

$$\log \left( \frac{p(x)}{1-p(x)} \right) = \beta x \quad \text{goal} \rightarrow \text{estimate } \hat{\beta}$$

For class 1: estimate  $\hat{\beta}$  such that  $p(\hat{\beta})$  is close to 1 as possible.

For class 0: estimate  $\hat{\beta}$  such that  $p(1-\hat{\beta})$  is close to 1 as possible.

$$L(\beta) = \prod_{y_i=1} p(x_i) \times \prod_{y_i=0} (1-p(x_i)) \quad l(\beta) = \sum_{i=1}^n y_i \log \left( \frac{1}{1+e^{-\beta x_i}} \right) + (1-y_i) \log \left( \frac{e^{-\beta x_i}}{1+e^{-\beta x_i}} \right)$$

Log likelihood function  $l(\beta)$  that needs to be optimized

$$l(\beta) = \sum_{i=1}^n y_i \beta x_i - \log(1+e^{\beta x_i}) \quad \beta = \arg \max l(\beta)$$

We consider Newton-Raphson to compute  $\beta$ . First two terms of Taylor series expansion.

$$\nabla_{\beta} l(\beta) = \nabla_{\beta} l(\beta^*) + (\beta - \beta^*) \nabla_{\beta\beta} l(\beta^*) \quad \nabla_{\beta} l(\beta^*) + (\beta - \beta^*) \nabla_{\beta\beta} l(\beta^*) = 0 \Rightarrow \beta^{t+1} = \beta^t - \frac{\nabla_{\beta} l(\beta^t)}{\nabla_{\beta\beta} l(\beta^t)} \rightarrow \text{Gradient w.r.t } \beta.$$

$$\boxed{\nabla_{\beta} l = \sum_{i=1}^n [y_i - p(x_i)] x_i}, \quad \nabla_{\beta\beta} l(\beta^t) = \nabla_{\beta} \cdot \sum_{i=1}^n [y_i - p(x_i)] x_i^T$$

$$\nabla_{\beta\beta} l(\beta^t) = - \sum_{i=1}^n p(x_i) [1 - p(x_i)] x_i x_i^T$$

$$\beta^{(t+1)} = \beta^{(t)} + C x^T w^{(t)} x^{-1} \times (Y - \hat{Y}^{(t)}) \quad \rightarrow \text{Repeat for } n \text{ times.}$$

$$\boxed{p(x) = \frac{1}{1+e^{-\beta x}}}$$

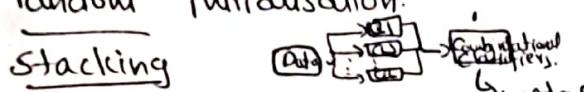
### Summary: Logistic regression

- Directly represent posterior distribution  $p(\phi | C_k)$
- Require fewer parameters than modelling the likelihood + prior

Ensemble - Use same training algorithm several times on subset of data.  
 well suited for unstable learning algorithms. Eg:- Decision trees, neural networks  
 stable learning algo - Nearest neighbour, linear regression.

5.8 5-fold cross-validation - Bagging - Bootstrap aggregation

$M=N$ , 63.2%, Perform multiple runs of the algo with different random initialisation.



- simplicity
- correlation b/w classifiers
- Feature combination

### Model combination

Eg: Mixture of models

Uncertainty which mixture component responsible for  $p(x) = \sum_{k=1}^K \pi_k N(x | \mu_k, \Sigma_k)$

Marginal probability of a data set  $p(x) = \prod_{n=1}^N \sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k)$

Bayesian model averaging: Suppose we have  $H$  different models  $h=1, \dots, H$   
 Marginal distribution  $p(x) = \sum_{h=1}^H p(x|h) p(h)$   $p(h)$  - prior probability

Model combination → Different data generated by different model components.

Uncertainty is about which component created which data points  
 one latent variable  $z_n$  for each data point.

$$p(x) = \prod_{n=1}^N p(x_n) = \prod_{n=1}^N \sum_{z_n} p(x_n | z_n)$$

### Bayesian model averaging

→ The whole data set is generated by a single model.

→ Uncertainty is about which model was responsible.

→ one latent variable  $z$  for the entire data set.

$$p(x) = \sum_z p(x, z)$$

$$y_{\text{com}}(x) = \frac{1}{M} \sum_{m=1}^M y_m(x) \quad y(x) = h(x) + \epsilon(x)^{\text{error}}$$

$$\text{Sum of square errors } E_x = [ \sum_m (y_m(x) - h(x))^2 ] = E_x [ \epsilon_m(x)^2 ]$$

$$\text{Average error individual models } E_{\text{AV}} = \frac{1}{M} \sum_{m=1}^M E_x [ \epsilon_m(x) ]^2$$

$$\text{Average error of committee } E_{\text{com}} = E_x \left[ \left\{ \frac{1}{M} \sum_{m=1}^M (y_m(x) - h(x)) \right\}^2 \right] = E_x \left[ \left\{ \frac{1}{M} \sum_{m=1}^M \epsilon_m(x) \right\}^2 \right]$$

$$E_{\text{com}} = \frac{1}{M} E_{\text{AV}} \text{ assumption: } E_x [ \epsilon_m(x) ] = 0 \rightarrow \text{zero mean}$$

$$E_x [ \epsilon_m(x) \epsilon_j(x) ] = 0$$

errors are uncorrelated

Ensemble method - simple method for improving classification

### Apparent contradiction

- classifiers should be trained on sample from the distribution on which it will be tested.

- Resampling seems to violate this recommendation.

### Explanation

- we do not attempt to model the full category distribution here.

- Try to find the decision boundary more directly.

- Increasing number of component classifiers broadens the class of implementable decision functions.

Boosting - combine multiple rules of thumb to make accurate decision.

Is problem learnable? - Probably Approximately Correct model (PAC)

chance of getting bad data? lower than threshold  $1-\delta$ ?

PAC - A problem is PAC learnable, if learning algorithm can find solution with some error  $\epsilon$  with probability greater than  $(1-\delta)$

Strong learner - error < threshold ( $\epsilon$ ), probability  $>1-\delta$ .

Ex:- sound recognition, lots of data, lots of parameters and hardware requirement, what if we don't have good hardware, or less dataset and how do we learn now?

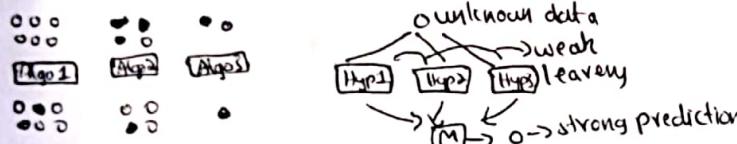
That's where weak learners comes into picture.

If a strong learner can solve a problem, a collection of weak learners can do it too,  
↳ This is based on hypothesis boosting mechanism.

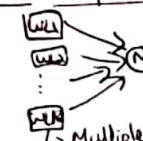
Hypothesis: Represent what model has learned.

For ex:- in supervised learning; final equation with weights  $y = b + w_1x_1 + \dots + w_px_p$

Instead we construct 3 hypotheses and finally do majority vote.



For complex problem



Initially all hypothesis have equal weightage to make a majority voting.

But what if each hypothesis added their importance changes or adapt depending on the errors of new hypothesis added. This is where adaptive boosting comes in. (Adaboost)

Normal hypothesis boosting - No weight associated with it.

Adaboost - weighted.

Step 1: Initialise weight of each sample to be the same.

Step 2: Construct the first hypothesis by training on these datapoints.

Step 3: Determine training error of these hypothesis and update the weights.

weight update - Increase weight of misclassified sample,  
- Decrease weight of correctly classified sample.

Update the weight of hypothesis itself.

First iteration - since first hypothesis weight=1, As iteration goes more hypothesis are added and similarly weights are updated.

If not enough hypothesis is added then partially boosted model won't have good performance to capture pattern. so boosting prone to underfitting.

### Adaboost

- 1) Initialise weights to all training points.  $w_i = \frac{1}{N}$  E.g: If I have 10k datapoint each datapoint have weight  $1/10k$
- 2) Calculate error rate for each weak classifier
- 3) Pick classifier or decision stump with the lowest error rate.  
↓  
Now voting compute voting power for the classifier  $\alpha = \frac{1}{2} \log \frac{1-\epsilon}{\epsilon}$   
(This gives the classifier its weak classifier strength in terms of how much is it more accurate as compared to the other weak classifier.)
- 4) Append the classifier into my ensemble classifier of  $f(x)$ . Then question is asked is  $f(x)$  good enough for making a strong prediction or have I gone around multiple rounds to create a very strong classifier. or are there any good classifier left that I can create. To check whether I have completed enough round is to check if the best classifier has an error rate of one by two. (1/2)

$$\epsilon = \sum w_i \text{ if and grow multiple decision stumps or small decision trees calculate error by each decision tree.}$$

- Once it is evaluated, and you feel that there are multiple more rounds that you have to conduct in term of creating more weak classifiers with some voting power. Then next step is to update the weights.

↓  
⑥ Update weights of each point where the previous classifier went wrong.

$$w_{\text{new}} = \begin{cases} \frac{w_{\text{old}}}{\alpha t_n} & \text{if point classified correctly (weight \downarrow)} \\ \frac{w_{\text{old}}}{\alpha t_n} & \text{if point classified wrongly (weight \uparrow)} \end{cases}$$

→ Go back to step 2.

- x — x — x —
1. Initialisation: Set  $w_n^{(1)} = \frac{1}{N}$  for  $n=1, \dots, N$
  2. For  $m=1, \dots, M$  iterations.

a) Train a new weak classifiers  $h_m(x)$  using the current weighting coefficients  $w^{(m)}$  by minimizing the weighted error function.

$$J_m = \sum_{n=1}^N \cdot w_n^{(m)} I(h_m(x) \neq t_n) \quad h_m(x) \rightarrow \text{weak classifier}$$

$$I(A) = \begin{cases} 1 & \text{if } A \text{ is true} \\ 0 & \text{else} \end{cases}$$

b) Estimate the weighted error of this classifier on  $X$ :

$$\epsilon_m = \frac{\sum_{n=1}^N w_n^{(m)} \cdot I(h_m(x_n) \neq t_n)}{\sum_{n=1}^N w_n^{(m)}}$$

c) Calculate a weighting coefficients for  $h_m(x)$ :  $\alpha_m = \ln \left( \frac{1 - \epsilon_m}{\epsilon_m} \right)$

d) Update the weighting coefficients:

$$w_n^{(m+1)} = w_n^{(m)} \exp \{ \alpha_m I(h_m(x_n) \neq t_n) \}$$

e) Make final prediction using the final combined model  $H(x) = \text{sign} \left\{ \sum_{m=1}^M \alpha_m h_m(x) \right\}$

Exponential error function  $E = \sum_{n=1}^N \exp \{ -t_n f_m(x_n) \}$

where  $f_m(x)$  is a classifier (linear combination of base classifier  $h_l(x)$ )

$$f_m(x) = \frac{1}{2} \sum_{l=1}^M \alpha_l h_l(x) \quad \text{Goal: Minimize } E \text{ wrt } \alpha_l \text{ and } h_l(x)$$

$$\frac{dE}{dh_m(x_n)} = 0$$

$$E = \sum_{n=1}^N w_n^{(m)} \exp \left\{ -\frac{1}{2} t_n \alpha_m h_m(x_n) \right\} \quad \begin{aligned} t_n h_m(x_n) &= +1 & \text{correctly classified points:} \\ t_n h_m(x_n) &= -1 \rightarrow \text{misclassified.} \end{aligned}$$

AdaBoost minimizes exponential error in a sequential fashion.

Limitations:

Original AdaBoost sensitive to misclassified training data points.

- Because of exponential error function • Improved by GentleBoost.

Two class classifier.

- Multiclass extensions available.

## Decision tree:

Each node specifies a test for some attribute.  
Each branch corresponds to a possible value of the attribute.

## Limitations of decision trees

- ⇒ often produce noisy or weak classifiers.
- ⇒ Do not generalize too well.
- ⇒ Training data fragmentation:
  - As tree progresses, splits
- ⇒ Overtraining and under training
  - Deep trees: fit the training data well, will not generalize well to new test data.
  - Shallow trees: not sufficiently refined.
- ⇒ Stability
  - Trees are sensitive to details of the training points.
  - If a single data point is only slightly shifted, a radically different tree may come out.
- ⇒ Expensive learning step.
  - due to costly selection of optimal split.

## CART framework.

- GQ → Binary or multi? → which property
- when node is leaf? → how tree can be pruned?
- How to deal with impure nodes → how missing attributes handled?

## Decision tree complexity

Data points  $\{x_1, \dots, x_N\}$  - dimensionality 0  
Storage:  $O(N)$ , Test runtime:  $O(\log N)$

Training runtime:  $O(DN^2 \log N)$

- Most expensive part.
- Critical step: selecting the optimal splitting point

## Best way to split the data.

- Don't look for global optimal split.
- Instead randomly use subset of K attributes on which to base the split.
- choose best splitting attribute.

E.g. maximizing the information gain

$$\Delta E = \sum_{k=1}^K \frac{|S_k|}{|S|} \cdot \sum_{j=s}^N p_j \log(p_j)$$

## Ensemble combination

combine output of several trees by averaging their posteriors

$$p(C|x) = \frac{1}{L} \sum_{l=1}^L p_{l,M_l}(C|x) \quad (l, M_l) \rightarrow \text{Tree leaves}$$

Bootstrap- subsampling with replacement.  
 Decision tree → Bagged decision tree → Random forest.

$\text{Ex: } \{1, 2, 3, \dots\} \cdot M_1, M_2, \dots, M_b = \frac{\sum_{i=1}^n M_i}{N}$  (N-samples)

why bootstrap? average over average reduces variance of the predicted mean.

Decreasing variance of prediction will reduce overfitting.

Decision tree → 1lp sample  
 categorical 1lp → classification tree  
 continuous 1lp → Regression trees  
 CART → classification and regression trees helps to construct predictive model of data.

### Bagged decision tree

Consider 1000 samples let's say we want to train B classifiers. For each of B classifier we split the dataset with replacement (Bootstrap). Then each classifier get 300 samples. Now construct B decision tree using corresponding samples. we pass the data to all tree and majority 1lp rules.

In case of regression, it takes the mean of all decision tree outputs.

why are we not pruning individual trees?  
 - won't they exhibit a high variance. Yes they will  
 we don't care because results are aggregated anyways. This aggregation lead to decrease in variance.

why use bagging in high variance ML algo?  
 High variance ML algo the B decision trees will have different structures and hence will predict differently. Hence aggregate prediction will decrease the variance of prediction.  
 However if we apply it to low variance ML algo like logistic regression then all B classifier look the same.

Suppose we have 1000 samples and we pick 600 of them then we left with 400 datapoints that decision tree haven't seen. These 400 samples are called out of bag samples (OOB).  
 OOB can be used to validate the random forest classifier.

### Bagging:

"Bootstrapped aggregation."  
 "Bootstrapped to high variance algorithm."  
 high variance- Different training samples leads to different results.

high variance model in decision trees  
 Different data fed into can change the structure of the data drastically.  
 we can reduce it with the help of bagging

### Random forest

The problem with bagged decision tree each B node can be split into on the same P features. This could lead to structural similarities b/w the trees. and hence correlation. However we want our model has to be different and correlated as possible.  
 Hence random forest are like bagged decision trees. but every tree can split based on subset of features M. this subset of features is different for each classifier. For classification tree with P features every tree should be able to split on  $\sqrt{P}$  features.

For regression trees each tree should be able to split on  $P/3$  features.

$$m = \begin{cases} \sqrt{P} & \text{classification} \\ P/3 & \text{regression} \end{cases}$$

**Standard perceptron**

**Multi-class perceptron**

**Non-linear basis function**

**Perception learning!**

- If olp is correct, leave the weights alone.
- If the olp unit incorrectly olp a zero, add the ilp vector to the weight vector.
- If the olp unit incorrectly olp a one, subtract the ilp vector from the weight vector.

$w_{kj}^{(t+1)} = w_{kj}^{(t)}$  This is guaranteed to converge to a correct solution if such solution exist.

$w_{kj}^{(t+1)} = w_{kj}^{(t)} - \eta(\text{err. } y_k(x_n; w) - t_{kn}) \phi_j(x_n)$  → Delta rule a.k.a LMS rule

**Loss function:**

- L2 Loss  $L(t, y(x)) = \sum_n (y(x_n) - t_n)^2$  Least square regression
- L1 loss  $L(t, y(x)) = \sum_n |y(x_n) - t_n|$  Median regression
- Cross-entropy loss  $L(t, y(x)) = -\sum_n \{t_n \ln y_n + (1-t_n) \ln(1-y_n)\}$
- Hinge loss  $L(t, y(x)) = -\sum_n [1 - t_n y(x_n)]$  → SVM classification.
- Softmax loss  $L(t, y(x)) = -\sum_n \sum_k \{t_k \ln \frac{\exp(y_k(x))}{\sum_j \exp(y_j(x))}\}$

**Regularization:**  $E(w) = \sum_n L(t_n, y(x_n; w)) + \lambda \|w\|^2$  → weight decay →  $\alpha$

**Avoid overfitting**

**Limitation of perception**

- Given right ilp feature perception can train any separable function.
- For some task exponential number of ilp feature is required.
- XOR problem.
- Non linear classifier based on the right kind of features, the problem become solvable.

**Multi-layer perceptron**

**Activation function:**  $g^{(0)}(a) = \sigma(a)$ ,  $g^{(l)}(a) = a$ .

Need an efficient way of adapting all weights, not just the last layer.  
Learning weights of hidden layers = learning features. → Gradient descent.

**Error function:**  $E(w) = \sum_n L(t_n, y(x_n; w)) + \lambda R(w)$

**G1** Obtain gradient for each weight

**G2** Adjust to new weight

Compute gradient for each weight

$$\frac{\partial E(w)}{\partial w_{ij}} = \sum_n \frac{\partial L(t_n, y(x_n; w))}{\partial w_{ij}} \rightarrow \text{L2-loss.}$$

$$R(w) = \|w\|_F^2 \rightarrow \text{L2-regulariser} \rightarrow \text{update weight } \frac{\partial E(w)}{\partial w_{ij}}$$

$\frac{\partial z}{\partial x} = \sum_{i=1}^L \frac{\partial z}{\partial y_i} \frac{\partial y_i}{\partial x}$

with increase in depth, there will be exponentially many paths, → infeasible to compute this way

**Numerical differentiation**

Given current state  $w^{(t)}$  make small change to  $w^{(t)}$  and except those that improve  $E(w^t)$

Horribly inefficient! Need several forward passes for each weight. Each forward pass is one run over the entire dataset.

**Naive analytical differentiation**

So we go for 3rd approach → Backpropagation.

## Gradient descent

### Two main steps

1. Computing the gradient for each weight
2. Adjusting the weights in the direction of the gradient

$$w_{kj}^{(t+1)} = w_{kj}^{(t)} - \eta \frac{\partial E(w)}{\partial w_{kj}} \Big|_{w(t)}$$

### Batch learning vs Stochastic learning.

#### Batch learning

- Process the full dataset at once to compute the gradient.

#### Stochastic learning

- Choose a single example from the training set
- Compute the gradient only based on this example
- This estimate will generally be noisy, which has some advantage.

#### Batch learning advantage

- Conditions of convergence are well understood.
- Many acceleration techniques (e.g. conjugate gradient) only operate in batch learning
- Theoretical analysis of the weight dynamics and convergence rates are simpler.

#### Stochastic learning advantage

- Usually much faster than batch learning
- often results in better solutions.
- Can be used for tracking changes.

### Minibatches

#### Idea

- Process only a small batch of training examples together.

- Start with a small batch size & increase it as training proceeds.

Advantage: • Gradient will more stable than for stochastic gradient descent, but still faster to compute than with batch learning.

- Take advantage of redundancies in the training set.

- Matrix operations are more efficient than vector operations.

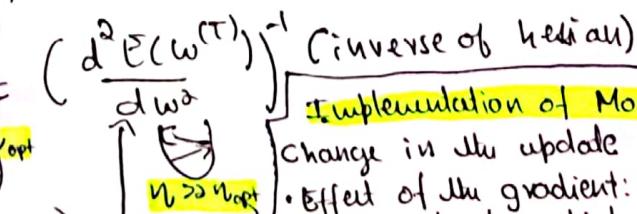
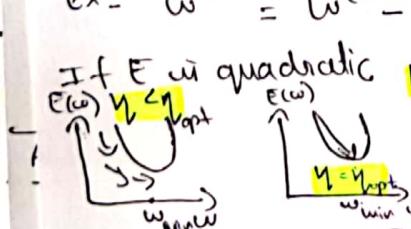
Caveats: • Error function should be normalised by the minibatch size, s.t. we can keep the same learning rate b/w minibatches.

$$E(w) = -\frac{1}{N} \sum_n L(t_n, y(x_n; w)) + \frac{\lambda}{N} R(w)$$

### Choose right learning rate

analysing the convergence of Gradient descent.

Ex:  $w^{(t+1)} = w^{(t)} - \eta \frac{dE(w)}{dw}$  what is the optimal  $\eta_{opt}$ ?



why learning is slow?

If the lips are correlated.

- The ellipse will be very elongated
- The direction of steepest descent is almost  $\perp$  to the direction towards minimum.



### The momentum method

- Idea: Instead of using the gradient to change the position of the weight "particle" use it to change velocity.

Intuition: Ball rolling on error surface.

- It starts off by following the error surface, but over time has accumulated momentum, it no longer does steepest descent.

#### Implementation of Momentum method:

Change in the update equations

- Effect of the gradient: increased the previous velocity, subject to a decay by  $\alpha$

$$v(t) = \alpha v(t-1) - \epsilon \frac{\partial E}{\partial w}(t)$$

Set weight change to the current velocity

$$\Delta w = v(t) = \alpha v(t-1) - \epsilon \frac{\partial E}{\partial w}(t)$$

$$= \alpha \Delta w(t-1) - \epsilon \frac{\partial E}{\partial w}(t)$$

#### Idea of RMSprop

- Divide the gradient by a running average of its recent magnitude.

$$\text{MeanSq}(w_{ij}; t) = 0.9 \text{MeanSq}(w_{ij}, t-1) +$$

$$\cdot \text{MeanSq}(w_{ij}, t) = 0.1 \left( \frac{\partial E}{\partial w_{ij}}(t) \right)^2$$

$$\cdot \text{MeanSq}(w_{ij}, t) = \frac{1}{10} \text{MeanSq}(w_{ij}, t) + 0.9 \left( \frac{\partial E}{\partial w_{ij}}(t) \right)^2$$

Backpropagation → idea: Compute the gradient layer by layer.

- each layer below builds upon the results of the layer above.
- the gradient is propagated backwards through the layers.

1) Convert the discrepancy b/w each output and its target value into an error derivative

$$E = \frac{1}{2} \sum_{j \in \text{outputs}} (t_j - y_j)^2$$

2) Compute error derivative in each hidden layer from error derivatives in the layers above.  $\frac{\partial E}{\partial y_i} = -(t_i - y_i)$

3) Use error derivatives w.r.t activities to get error derivatives w.r.t the incoming weights.

1  $y_j \rightarrow \text{old of layer } j$       convolution  
 $z_j \rightarrow \text{inp of layer } j$        $z_j = \sum_i w_{ij} y_i$   
 $y_i = g(z_j)$

$$\frac{\partial E}{\partial y_i} \rightarrow \frac{\partial E}{\partial w_{ik}}$$

### Forward pass

$y^{(0)} = x$   
for  $k=1, \dots, l$  do  
 $z^{(k)} = w^{(k)} y^{(k-1)}$   
 $y^{(k)} = g_k(z^{(k)})$   
end for  
 $y = y^{(l)}$   
 $E = L(t, y) + \lambda R(w)$

### Backward pass

$h \leftarrow \frac{\partial E}{\partial y} = \frac{\partial}{\partial y} [L(t, y) + \lambda R(w)]$

for  $k=l, l-1, \dots, 1$  do

$$h \leftarrow \frac{\partial E}{\partial z^{(k)}} = h \odot g'(y^{(k)})$$

$$\frac{\partial E}{\partial w^{(k)}} = h y^{(k-1)T} + \lambda \frac{\partial E}{\partial w^{(k)}}$$

$$h \leftarrow \frac{\partial E}{\partial y^{(k-1)}} = w^{(k)T} h$$

end for

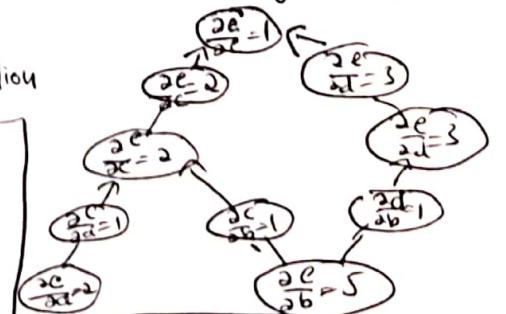
Forward mode differentiation ( $\frac{\partial}{\partial x}$ )

Apply operator  $\frac{\partial}{\partial x}$  to every node

Reverse mode differentiation ( $\frac{\partial z}{\partial x}$ )

why do we care?

- Let's consider the example again
- Using reverse-mode differentiation from e up ...
  - Runtime:  $O(\#edges)$
  - Result: derivative of e with respect to every node.



### Automatic differentiation:

softmax output

$$E(w) = - \sum_{n=1}^N \sum_{k=1}^K \{ I(t_n=k) \ln \frac{\exp(w_k^T x)}{\sum_{j=1}^K \exp(w_j^T x)} \}$$

### Practical issues

- Exponentials get very big and can have vastly different magnitude
- Trick 1: Do not compute first softmax, then log, but instead directly evaluate  $\log \exp$  in the denominator.
- Trick 2: Softmax has the property that for a fixed vector b  $\text{softmax}(a+b) = \text{softmax}(a)$

RL: An agent by interacting with the environment and obtaining a reward signal for its actions.

Cart Pole balancing, context-aware object hand-over, opening doors.

Supervised learning: given input  $x$  and labels/output  $y$ . Learn: a mapping  $f(x) = y$ .

Unsupervised learning: input  $x$ , same Learn: some structure about the data

Reinforcement learning: Given: a set of experiences Learn: how to behave so that the reward is maximised.

A reinforcement modelled as Markov decision Process.

MDPC  $(S, A, T, R)$ , return  $R$ ,  $J = E \left[ \sum_t r(s_t, a_t) \right]$

Policy  $\pi: S \rightarrow A$ , parameter  $\pi_\theta$ , trajectory  $T = (s_0, a_0, s_1, \dots, a_n, s_{n+1})$

$P_\pi(T) = P(s_0) \prod_{i=0}^{n-1} P_\pi(a_i | s_i) P_{T+1}(s_{i+1} | a_i)$  (episode or rollout)

Stochastic policy  $\rightarrow$  Sampling from distribution of action state  $a_t \sim \pi(s_t)$

Deterministic policy  $\rightarrow$  Fixed,  $a_t = \pi(s_t)$  e.g: Grasping a door. Eg: stroke in tennis

Finite & infinite horizon tasks  $\rightarrow$  finite tasks periodic ends after a certain number of time steps or when a goal is reached. Eg: Navigation

$$R(T) = \sum_{t=0}^T r(s_t, a_t)$$

$$\text{Return, } R(T) = \sum_{t=0}^T r(s_t, a_t)$$

discontinued by a factor  $\gamma \in [0, 1]$

RL objective  $\rightarrow$  Maximize the agent's expected return  $\pi^* = \arg \max_{\pi} E \left[ \sum_t r(s_t, a_t) \right]$

$\downarrow$  find  $\pi^*$   $\rightarrow$  Maximize the parameterised policy

$$\pi^* = \arg \max_{\theta} E_{\pi \sim \pi_\theta} \left[ \sum_t r(s_t, a_t) \right] = \arg \max_{\theta} \int P(T|\pi) R(T) dT$$

Transition function  $T$  + reward function  $r$  = model

$\rightarrow$  If  $T$  and  $r$  are known in advance, model-based RL (Exploit the model knowledge)

$\rightarrow$  If  $T$  and  $r$  are unknown, model-free RL (Learn using trial error)

We often prefer combination of both

Exploration: Acting by trying out action that may not be optimal one.

Exploitation: Acting according to the best policy.

Tradeoff: if the agent exploits too much, it risks of converging to sub-optimal policy. (may not see entire environment)

if the agent does exploration, it may not know when to stop or converge. (Randomly choose policy) (slow)

Off-policy: choose policy different from what it has learned. (Randomly choose policy) (slow)

On-policy: an agent samples actions from the policy that it is trying to learn. (fast)

Value function: It tells if it's good to be in this state and not the other

State value function:  $V^\pi(s) \Rightarrow$  It is the expected return if we start in state  $s_0$  and then act according to the policy.  $V^\pi(s) = E[R(T) | s_0=s]$

State-action value function:  $Q^\pi(s, a)$  evaluates both state and action (Q-function)

$\Rightarrow$  Expected return for agent starting at state  $s_0$  and action  $a$  and act subsequently.

$$Q^\pi(s, a) = E_{T \sim \pi} [R(T) | s_0=s, a_0=a]. \quad V^\pi(s) = \arg \max_a Q^\pi(s, a) \quad \forall s \in S$$

Optimal state value function (w.r.t policy  $\pi$ ):  $V^*(s) = \arg \max_{\pi} V^\pi(s) \quad \forall s \in S, \forall a \in A(s)$

Optimal state-action value function:  $Q^*(s, a) = \arg \max_{\pi} Q^\pi(s, a) \quad \forall s \in S, \forall a \in A(s)$

$V^*(s) = \max_{a \in A(s)} Q^*(s, a) \quad a^* = \arg \max_a Q^*(s, a) \Rightarrow$  consequence.

By having optimal state-action we know how good each state is and how good a policy is.

Advantage function:  $A(s, a) = Q(s, a) - V(s) \Rightarrow$  Tells how good specific action is.

Used to reduce variance of exploration.

Model based learning: Aim: find optimal value function

$$V^*(s) = \max_{\pi} E \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right]$$

We can find optimal value function by solving Bellman equations.

$$V^*(s) = \max_a [r(s, a) + \gamma \sum_{s'} P(s'|s, a) V^*(s')] \quad (1)$$

We need to find how good state  $s'$  is. First we need immediate reward. Then we have discounted expectation of how good the subsequent states are. (2)

Optimal policy is then  $\pi^*(s) = \arg \max_a (r(s, a) + \gamma \sum_{s'} P(s'|s, a) V^*(s'))$

Two algorithms: 1) Value iterations 2) Policy iterations. (solving algorithm for optimal value function)

Value iteration: uses directly bellman equation and solves for them. (iterative algorithm). start with random initialisation of  $V(s)$ ,  $\forall s \in S$ , then iterate over the equations. Every state we maximize over the action. solve until convergence, Eg:-  $\forall s \in S, |V^i(s) - V^{i-1}(s)| \leq \epsilon$  for small  $\epsilon$ .

It is guaranteed to converge to an optimal policy.

Policy iteration: Here we modify the policy directly. so it is called policy iteration.

Here we randomly initialise a policy  $\pi$  and then solve the system of equations.

$$V^{\pi_i}(s) = r(s, \pi_i(s)) + \gamma \sum_{s'} P(s'|s, a) V^{\pi_i}(s')$$

$$\text{Update the policy } \pi_{i+1}(s) = \arg \max_a (r(s, a) + \gamma \sum_{s'} P(s'|s, a) V^{\pi_i}(s'))$$

Model free learning:

Setup: while learning collect trajectories  $T$  and the accompanying rewards.

use this setup to find optimal policy.

Two types of algorithms: Temporal difference learning; perform value/policy update at every step.

Monte carlo learning: estimate the return and then perform value/policy updates.

Temporal difference - TD( $\lambda$ )-Learning: In TD( $\lambda$ ) we try to bring value function closer to the reward function.

The parameter  $\lambda$  controls the amount of prediction during learning.

$$V(s_t) = V(s_t) + \alpha \cdot C(r(s_t, a_t) + \gamma V(s_{t+1}) - V(s_t)), \text{ where } \alpha \text{ is a learning rate.}$$

Q learning: Here we try to estimate the state-action value function.

$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha (r(s_t, a_t) + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t))$$

Policy search (PS)

Here we don't need to find value function explicitly rather we would optimise policy space directly. Useful for parameterised policies. They help to incorporate prior knowledge about the policy.

Three type of PS algorithms:

Policy gradient method

We represent policy as differentiable function then we use gradient based method to optimise.

Information-theoretic algo

We want to do policy update without destroying the previous policies. We try to bound the policy updates.

Expectation-maximisation-based method

Frame policy search as expectation maximisation and then solve EM.

## Policy gradients

Modify policy parameter  $\theta$  by calculating the gradient of the expected result.  $\theta \leftarrow \theta + \nabla J(\theta)$   $J(\theta) = \int R(\theta) P(T|\theta) d\theta$

Policy gradient algorithms make use of likelihood ratio trick.

$$\nabla_{\theta} P(T|\theta) = P(T|\theta) \nabla_{\theta} \log P(T|\theta)$$

## Reinforce algorithm

Initialise policy parameter  $\theta$  randomly

for  $i \leftarrow 1$  to  $N$  do (let agent explore for  $N$  iteration)  
 $T \rightarrow \{\}$  collect trajectories

for  $j \leftarrow 1$  to  $M$  do  
 $T \leftarrow \text{sample}(T|\theta)$  sample  $M$  trajectories  
 $\gamma \leftarrow \gamma \cup T \rightarrow \text{append}$

$$J_{\theta} \leftarrow \frac{1}{M} \sum_{j=1}^M \sum_t \gamma_t$$

$$\theta \leftarrow \theta + \nabla J_{\theta}$$

on  $i$ th iteration we estimate expected return  
 collect  $M$  trajectories.  
 estimate expected return.  
 estimate gradient of return.  
 update policy parameters.  
 repeat for  $N$  iteration

## Actor-Critic Learning

Value based algorithms - critic based

Policy search algorithms - actor based

we combine both of them together to have actor-critic family of RL.  
 which estimate the value function and maintain a policy at the same time.

Actor-critic algorithms make use of baseline  $b$  when estimating the gradient of  $J$ .

$$\nabla_{\theta} J(\theta) = E \left[ \sum_{t=0}^b \nabla_{\theta} \log P(a|s_t) (R_t - b_t) \right]$$

why do we do baseline  $\rightarrow$  we usually observe high variance of policy  
 baseline try to reduce  $b$  variance.