

## 1 Linear Classifier

Feature vectors  $x$ , labels  $y$

$$x \in \mathbb{R}^d$$
$$y \in \{-1, 1\}$$

Training set

$$S_n = \{(x^{(i)}, y^{(i)}), i = 1, \dots, n\}$$

Classifier

$$h: \mathbb{R}^d \rightarrow \{-1, 1\}$$
$$\chi^+ = \{x \in \mathbb{R}^d : h(x) = 1\}$$
$$\chi^- = \{x \in \mathbb{R}^d : h(x) = -1\}$$

Training error

$$\varepsilon_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{h(x^{(i)}) \neq y^{(i)}\}$$

Test error (over disjoint set of examples)

$$\varepsilon(h)$$

Set of classifiers

$$h \in H$$

### 1.1 Linear classifiers through origin

Set of all points that satisfies a line through the origin.

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$
$$X = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

Decision Boundary

$$\{x: \theta_1 x_1 + \theta_2 x_2 = 0\}$$
$$\{x: \theta \cdot X = 0\}$$

Linear Classifier through origin

$$h(x, \theta) = \text{sign}(\theta \cdot X)$$
$$\Theta \in \mathbb{R}^d$$

### 1.2 Linear classifiers

General linear Classifier (with Intercept)

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$
$$X = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

Decision Boundary

$$\{x: \theta \cdot X + \theta_0 = 0\}$$

Linear Classifier through origin

$$h(x, \theta, \theta_0) = \text{sign}(\theta \cdot X + \theta_0)$$
$$\theta \in \mathbb{R}^d$$
$$\theta_0 \in \mathbb{R}$$

### 1.3 Linear Separation

Traning examples  $S_n = \{(x^{(i)}, y^{(i)}), i = 1, \dots, n\}$  are linear separable if there exists a parameter vector  $\theta$  and offset parameter  $\theta_0$  such that  $y^{(i)}(\theta \cdot x^{(i)} + \theta_0) > 0$  for all  $i = 1, \dots, n$ .

$$(\hat{\theta} \cdot x^{(i)}) > 0 \begin{cases} y^{(i)} > 0 \text{ and } \theta \cdot x^{(i)} > 0 \\ y^{(i)} < 0 \text{ and } \theta \cdot x^{(i)} < 0 \end{cases}$$

$y^{(i)}(\theta \cdot x^{(i)}) > 0$  if label and classified result match. This leads to a new definition of the **Training error**:

$$\varepsilon_n(\theta) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{y^{(i)}(\theta \cdot x^{(i)}) \leq 0\}$$
$$\varepsilon_n(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{y^{(i)}(\theta \cdot x^{(i)} + \theta_0) \leq 0\}$$

### 2 Perceptron through origin

Perceptron( $\{(x^{(i)}, y^{(i)}), i = 1, \dots, n\}, T$ ):

initialize  $\theta = 0$  (vector);

for  $t = 1, \dots, T$  do

for  $i = 1, \dots, n$  do

if  $y^{(i)}(\theta \cdot x^{(i)}) \leq 0$  then

update  $\theta = \theta + y^{(i)} x^{(i)}$

### 1.5 Perceptron with Offset

Perceptron( $\{(x^{(i)}, y^{(i)}), i = 1, \dots, n\}, T$ ):

initialize  $\theta = 0$  (vector);  $\theta_0 = 0$  (scalar)

for  $t = 1, \dots, T$  do

for  $i = 1, \dots, n$  do

if  $y^{(i)}(\theta \cdot x^{(i)} + \theta_0) \leq 0$  then

update  $\theta = \theta + y^{(i)} x^{(i)}$

update  $\theta_0 = \theta_0 + y^{(i)}$

### 1.6 Margin Boundary

The Margin Boundary is the set of points  $x$  which satisfy  $\theta \cdot x + \theta_0 = \pm 1$ . So, the signed distance from the decision boundary to the margin boundary is  $\frac{1}{\|\theta\|}$ .

$$\frac{y^{(i)}(\theta \cdot x^{(i)} + \theta_0)}{\|\theta\|} = \frac{1}{\|\theta\|}$$

Hinge Loss (agreement)

$$\text{Agreement} = z = y^{(i)}(\theta \cdot x^{(i)} + \theta_0)$$

$$\text{Loss}_h(z) = \begin{cases} 0 & \text{if } z \geq 1 \\ 1 - z & \text{if } z < 1 \end{cases}$$

**Regularization** means pushing out the margin boundaries by adding  $\max(\frac{1}{\|\theta\|})$  or  $\min(\frac{1}{2} \|\theta\|^2)$  to the objective function.

**Objective Function**

Objective function = average loss + regularization

Objective function is minimized, learning becomes an optimization problem. Using hinge loss and margin boundaries is called **Support Vector Machine** or **Large margin linear classification**:

$$J(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n \text{Loss}_h(z) + \frac{\lambda}{2} \|\theta\|^2$$

Where  $\lambda > 0$  is called the regularization parameter that regulates how important the margin boundaries are in comparison to the average hinge loss.

### 1.7 Gradient Descent

Assume  $\theta \in \mathbb{R}$  the goal is to find  $\theta$  that minimizes  $J(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n \text{Loss}_h(y^{(i)}(\theta \cdot x^{(i)} + \theta_0)) + \frac{\lambda}{2} \|\theta\|^2$  through gradient descent.

In other words, we will

- Start  $\theta$  at an arbitrary location:  $\theta \leftarrow \theta_{start}$
- Update  $\theta$  repeatedly with  $\theta \leftarrow \theta - \eta \frac{\partial J(\theta, \theta_0)}{\partial \theta}$  until  $\theta$  does not change significantly.

Where  $\eta > 0$  is called the stepsize or **learning parameter**.

### 1.8 Stochastic Gradient Descent

$$J(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n \text{Loss}_h(z) + \frac{\lambda}{2} \|\theta\|^2$$
$$= \frac{1}{n} \sum_{i=1}^n \left[ \text{Loss}_h(z) + \frac{\lambda}{2} \|\theta\|^2 \right]$$

With stochastic gradient descent, we choose  $i \in \{1, \dots, n\}$  at random and update  $\theta$  such that

$$\theta \leftarrow \theta - \eta \nabla_{\theta} \left[ \text{Loss}_h(z) + \frac{\lambda}{2} \|\theta\|^2 \right]$$

## 2 Linear Algebra

### 2.1 Distance

Consider a line  $L$  in  $\mathbb{R}^2$  given by the equation  $L: \theta \cdot x + \theta_0 = 0$  where  $\theta$  is a vector normal to the line  $L$ . Let the point  $P$  be the endpoint of a vector  $x_0$  (so the coordinates of  $P$  equal the components of  $x_0$ ).

The shortest distance  $d$  between the line  $L$  and the point  $P$  is:

$$d = \frac{|\theta \cdot x_0 + \theta_0|}{\|\theta\|}$$

## 3 Regression and Classification

Classification:

$$S_n = \{(x^{(t)}, y^{(t)}) | t = 1, \dots, n\}$$
$$x^{(t)} \in \mathbb{R}^d, y^{(t)} \in \{-1, 1\}$$

Regression:

$$y^{(t)} \in \mathbb{R}$$
$$f(x, \theta, \theta_0) = \sum_{i=1}^d \theta_i x_i + \theta_0 =$$
$$= \theta \cdot x + \theta_0$$

### 3.1 Objective for linear regression

The empirical risk  $R_n$  is defined as

$$R_n(\theta) = \frac{1}{n} \sum_{t=1}^n \text{Loss}(y^{(t)} - \theta \cdot x^{(t)})$$

where  $(x^{(t)}, y^{(t)})$  is the  $t$ th training example (and there are  $n$  in total), and Loss is some loss function, such as hinge loss. Possible to get **closed form solution** for gradient because function is concave. Only possible if the  $d \times d$  matrix  $A$  is invertible. Computationally expensive if dimensions are very high like in bag of words approach.

$$\nabla R_n(\theta) = A\theta - b (= 0)$$
$$= A^{-1} b$$

where

$$A = \frac{1}{n} \sum_{t=1}^n x^{(t)} (x^{(t)})^T$$
$$b = \frac{1}{n} \sum_{t=1}^n y^{(t)} x^{(t)}$$

$b$  is a vector with dimensionality  $d$ .

### 3.2 Gradient based Approach

Nudge gradient in the opposite direction to find (local) minima.

$$\nabla_{\theta} (y^{(t)} - \theta x^{(t)})^2 / 2 =$$
$$= (y^{(t)} - \theta x^{(t)}) \nabla_{\theta} (y^{(t)} - \theta x^{(t)}) =$$
$$= -(y^{(t)} - \theta x^{(t)}) \cdot x^{(t)}$$

- initialize  $\theta = 0$
- randomly pick  $t = \{1, \dots, n\}$
- $\theta = \theta + \eta (y^{(t)} - \theta x^{(t)}) \cdot x^{(t)}$

Where  $\eta$  is the learning rate (steps) and the learning rate gets smaller the closer you get  $\eta_k = \frac{1}{1+k}$

### 3.3 Ridge Regression

Regularization is trying to push away from perfect fit.

$$J_{n,\lambda}(\theta, \theta_0) = \frac{1}{n} \sum_{t=1}^n \frac{(y^{(t)} - \theta x^{(t)} - \theta_0)^2}{2} + \frac{\lambda}{2} \|\theta\|^2$$
$$\nabla_{\theta} (J_{n,\lambda}) = \lambda \theta - (y^{(t)} - \theta x^{(t)}) x^{(t)}$$

- initialize  $\theta = 0$
- randomly pick  $t = \{1, \dots, n\}$
- $\theta = \theta + \eta \lambda \theta - (y^{(t)} - \theta x^{(t)}) x^{(t)} = (1 - \eta \lambda) \theta + \eta (y^{(t)} - \theta \cdot x^{(t)})$

## 3.4 Kernels

$$\phi(x) = [x_1, x_2, x_1^2, \sqrt{2} x_1 x_2, x_2^2]$$
$$\phi(x') = [x'_1, x'_2, x'^2_1, \sqrt{2} x'_1 x'_2, x'^2_2]$$
$$\phi(x) \cdot \phi(x') = x_1 x'_1 + x_2 x'_2 + x_1 x'^2_1 + 2 x_1 x'_1 x_2 x'_2 + x_2 x'^2_2$$
$$= (x_1 x'_1 + x_2 x'_2) + (x_1 x'_1 + x_2 x'_2)^2$$
$$= (x_1 x'_1 + x_2 x'_2) + (x_1 x'_1 + x_2 x'_2)^2$$

### 3.5 Kernel Perceptron

The parameter vector of a perceptron algorithm can also be written as:

$$\theta = \sum_{j=1}^n \alpha_j y^{(j)} \phi(x^{(j)})$$

Where  $\alpha_j$  represents the number of classification mistakes the perceptron algo made. Every time a misclassification happens the parameter vector is updated with the product of the label and the feature vector  $\theta = \theta + y^{(i)} \phi(x^{(i)})$ . The goal of the Kernel Perceptron algo is to find the vector  $\alpha$  with the counts of the misclassifications.

**Kernel Perceptron**( $\{(x^{(i)}, y^{(i)}), i = 1, \dots, n, T\}$ )

initialize  $\alpha_1, \alpha_2, \dots, \alpha_n$ ; to some values

for  $t = 1, \dots, T$  do

for  $i = 1, \dots, n$  do

if  $y^{(i)} \sum_{j=1}^n \alpha_j y^{(j)} K(x^{(i)}, x^{(j)}) \leq 0$  then

update  $\alpha_i = \alpha_i + y^{(i)}$

## 3.6 Radial basis Kernel

$$K(x, x') = e^{-\frac{1}{2} \|x - x'\|^2}$$

## 4 Recommender Systems

### 4.1 K nearest neighbors

The  $K$ -Nearest Neighbor method makes use of ratings by  $K$  other "similar users when predicting  $Y_{ai}$ .

Let  $KNN(a)$  be the set of  $K$  users "similar to user  $a$ , and let  $\text{sim}(a, b)$  be a similarity measure between users  $a$  and  $b \in KNN(a)$ . The  $K$ -Nearest Neighbor method predicts a ranking  $Y_{ai}$  to be:

$$\hat{Y}_{ai} = \frac{\sum_{b \in KNN(a)} \text{sim}(a, b) Y_{bi}}{\sum_{b \in KNN(a)} \text{sim}(a, b)}$$

The similarity measure  $\text{sim}(a, b)$  could be any distance function between the feature vectors  $x_a$  and  $x_b$  of users  $a$  and  $b$ , e.g. the euclidean distance  $\|x_a - x_b\|$  and the cosine similarity  $\cos \theta = \frac{x_a \cdot x_b}{\|x_a\| \|x_b\|}$ .

### 4.2 Collaborative Filtering

Matrix  $Y$  with  $n$  rows (users) and  $m$  columns (Movies) is sparse (entries missing),  $(a, i)$ th entry  $Y_{ai}$  is the rating by user  $a$  of movie  $i$  if this rating has already been given, and blank if not. Goal is to predict matrix  $X$  with no missing entries.

Let  $D$  be the set of all  $(a, i)$ 's for which a user rating  $Y_{ai}$  exists, i.e.  $(a, i) \in D$  if and only if the rating of user  $a$  to movie  $i$  exists.

$$J = \sum_{(a,i) \in D} \frac{(Y_{ai} - \frac{1}{2} UV^T)_{ai}^2}{2} + \frac{\lambda}{2} \left( \sum_{a,k} U_{ak}^2 + \sum_{i,k} V_{ik}^2 \right)$$
$$u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}; v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$
$$v = \begin{bmatrix} 2 \\ 7 \\ 8 \end{bmatrix}$$
$$uv^T = \begin{bmatrix} 2u_1 7u_1 8u_1 \\ 2u_2 7u_2 8u_2 \end{bmatrix}$$

Take derivative of Objective function  $J$  with respect to every user, set it to zero and find respective  $u_i$  value:

$$\frac{d}{du_1} \left( \frac{(7-8u_1)^2}{2} + \frac{\lambda}{2} u^2 \right) = 0$$
$$\frac{d}{du_2} (J) = 0$$
$$u_1 = \frac{66}{\lambda+68}; u_2 = \frac{16}{\lambda+53}$$

Use resulting values for  $u$  to compute  $uv^T$  compare resulting matrix  $X$  with matrix  $Y$  and start again. Continue until convergence.

## 5 Clustering

Two Views:

Clustering input:  $S_n = \{x^{(i)} | i = 1, \dots, n\}$

Clustering output are indexes for the data that partition the data:  $C_1, \dots, C_k$ ; where  $C_1 \cup C_2 \cup \dots \cup C_k = \{1, 2, \dots, n\}$  and the union of all  $C_j$ 's is the original set and the intersection of any  $C_i$  and  $C_j$  is an empty set.

Representatives of clusters:  $z^{(1)}, \dots, z^{(k)}$ .

Cost of partitioning is the sum of costs of individual clusters:  $\text{cost}(C_1, \dots, C_k) = \sum_{j=1}^k \text{cost}(C_j)$ .

Cost of cluster is sum of distances from data points to the representative of the cluster:  $\text{Cost}(C, z) = \sum_{i \in C} \text{distance}(x^{(i)}, z)$

Cosine similarity:  $\cos(x^{(i)}, x^{(j)}) = \frac{x^{(i)} \cdot x^{(j)}}{\|x^{(i)}\| \|x^{(j)}\|}$  is not sensitive to magnitude of vector (will not react to length).

Euclidean square distance:  $\text{dist}(x^{(i)}, x^{(j)}) = \|x^{(i)} - x^{(j)}\|^2$ . Will react to length.

$$\text{cost}(C_1, \dots, C_k; Z^{(1)}, \dots, Z^{(1)}) = \sum_{j=1}^k \sum_{i \in C_j} \|x^{(i)} - z^{(j)}\|^2$$

### 5.1 The K-Means Algorithm

Only works with Euclidean square distance.

Given a set of feature vectors  $S_n = \{x^{(i)} | i = 1, \dots, n\}$  and the number of clusters  $K$  we can find cluster assignments  $C_1, \dots, C_K$  and the representatives of each of the  $K$  clusters  $z_1, \dots, z_K$ :

1. Randomly select  $z_1, \dots, z_K$
2. Iterate
  - (a) Given  $z_1, \dots, z_K$ , assign each data point  $x^{(i)}$  to the closest  $z_j$ , so that  $\text{Cost}(z_1, \dots, z_K) = \sum_{i=1}^n \min_{j=1, \dots, K} \|x^{(i)} - z_j\|^2$
  - (b) Given  $C_1, \dots, C_K$  find the best representatives  $z_1, \dots, z_K$ , i.e. find  $z_1, \dots, z_K$  such that  $z_j = \arg\min_z \sum_{i \in C_j} \|x^{(i)} - z\|^2$

The best representative is found by optimization (gradient with respect to  $z^{(j)}$ , setting to zero and solving for  $z^{(j)}$ ). It is the centroid of the cluster:  $z^{(j)} = \frac{\sum_{i \in C_j} x^{(i)}}{|C_j|}$

The clustering output that the K-Means algorithm converges to depends on the initialization.

### 5.2 K-Medoids Algorithm

Finds the cost-minimizing representatives  $z_1, \dots, z_K$  for any distance measure. Uses real data points for initialization.

1. Randomly select  $\{z_1, \dots, z_K\} \subseteq \{x_1, \dots, x_n\}$
2. Iterate
  - (a) Given  $z_1, \dots, z_K$ , assign each data point  $x^{(i)}$  to the closest  $z_j$ , so that  $\text{Cost}(z_1, \dots, z_K) = \sum_{i=1}^n \min_{j=1, \dots, K} \|x^{(i)} - z_j\|^2$
  - (b) Given  $C_j \in \{C_1, \dots, C_K\}$  find the best representative  $z_j \in \{x_1, \dots, x_n\}$  such that  $\sum_{x^{(i)} \in C_j} \text{dist}(x^{(i)}, z_j)$  is minimal

### 6 Generative Models

Understand structure of data probabilistically.

#### 6.1 Multinomial Models

Fixed Vocabulary  $W$

Multinomial model  $M$  to generate text in documents.

Document  $D$

Likelihood of generating certain word  $w \in W$ :  $p(w|\theta) = \theta_w$  where  $\theta_w \geq 0$  and  $\sum_{w \in W} \theta_w = 1$ .

Likelihood function:

$$P(D|\theta) = \prod_{i=1}^n \theta_{w_i} = \prod_{w \in W} \theta_w^{\text{count}(w)}$$

Toy Example:

$$\begin{aligned} \theta_1 : \theta_{cat} &= 0.3; \theta_{dog} = 0.7 \\ \theta_2 : \theta_{cat} &= 0.9; \theta_{dog} = 0.1 \\ D &= \{cat, cat, dog\} \\ P(D|\theta_1) &= 0.3^2 \cdot 0.7 = 0.063 \\ P(D|\theta_2) &= 0.9^2 \cdot 0.1 = 0.081 \end{aligned}$$

Maximum likelihood:

$$\begin{aligned} \max_{\theta} P(D|\theta) &= \max_{\theta} \prod_{w \in W} \theta_w^{\text{count}(w)} \\ \log \prod_{i=1}^n \theta_w^{\text{count}(w)} &= \sum_{w \in W} \text{count}(w) \log(\theta_w) \\ W &= \{0, 1\}; \theta_0 = \theta; \theta_1 = (1 - \theta) \\ \frac{d}{d\theta} (\text{count}(0) \log(\theta) + \text{count}(1) \log(1 - \theta)) &= \\ &= \frac{\text{count}(0)}{\theta} - \frac{\text{count}(1)}{1 - \theta} = 0 \\ \hat{\theta} &= \frac{\text{count}(0)}{\text{count}(0) + \text{count}(1)} \end{aligned}$$

For any length of  $W$ :

$$\hat{\theta} = \frac{\text{count}(w)}{\sum_{w' \in W} \text{count}(w')}$$

### 6.2 Prediction

Goal: categorize between minus and plus class. Both classes have a associated parameter  $\theta^+$  and  $\theta^-$   
Class conditional distribution:

$$\log\left(\frac{P(D|\theta^+)}{P(D|\theta^-)}\right) = \begin{cases} \geq 0, & + \\ < 0, & - \end{cases}$$

Model is the same as a linear classifier through origin:

$$\begin{aligned} \log(P(D|\theta^+)) - \log(P(D|\theta^-)) &= \\ &= \log \prod_{w \in W} \theta_w^{+\text{count}(w)} - \log \prod_{w \in W} \theta_w^{-\text{count}(w)} = \\ &= \sum_{w \in W} \text{count}(w) \log(\theta_w^{+\text{count}(w)}) - \sum_{w \in W} \text{count}(w) \log(\theta_w^{-\text{count}(w)}) = \\ &= \sum_{w \in W} \text{count}(w) \log \frac{\theta_w^{+\text{count}(w)}}{\theta_w^{-\text{count}(w)}} \\ &= \sum_{w \in W} \text{count}(w) \tilde{\theta}_w \end{aligned}$$

### 6.3 Prior, Posterior and Likelihood

From bayes rule  $P(A|B) = \frac{P(B|A)P(A)}{P(B|B)}$  we get:

$$P(y = +|D) = \frac{P(D|\theta^+)P(y=+)}{P(D)}$$

Where  $P(y = +|D)$  is the posterior distribution and  $P(y = +)$  is the prior distribution while  $P(D|\theta^+)$  is the likelihood of document  $D$  given parameter  $\theta^+$ . This yields (after some work) a linear separator with offset:

$$\begin{aligned} \log\left(\frac{P(y=+|D)}{P(y=-|D)}\right) &= \frac{P(D|\theta^+)P(y=+)}{P(D|\theta^-)P(y=-)} \\ &= \log\left(\frac{P(D|\theta^+)}{P(D|\theta^-)}\right) + \log\left(\frac{P(y=+)}{P(y=-)}\right) \\ &= \sum_{w \in W} \text{count}(w) \tilde{\theta}_w + \tilde{\theta}_0 \end{aligned}$$

### 6.4 Gaussian Generative models

Vectors in  $x \in \mathbb{R}^d$  "cloud" data in which  $\mu$  (average over all points) is the center of the cloud and  $\sigma^2$  (square of average distance) the radius.  
Probability of  $x$  generated by gaussian cloud:

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

Likelihood of the training data:  $S_n = \{x^{(i)} | i = 1, \dots, n\}$  given the gaussian model  $p(S_n|\mu, \sigma^2) = \prod_{i=1}^n P(x^{(i)}|\mu, \sigma^2)$ .  
To get the MLE calculate likelihood, take the log and massage:

$$\begin{aligned} \log\left(\prod_{i=1}^n \frac{1}{(2\pi\sigma^2)^{d/2}} \exp\left(-\frac{1}{2\sigma^2} \|x - \mu\|^2\right)\right) &= \\ &= \sum_{i=1}^n \log \frac{1}{2\sigma^2} + \sum_{i=1}^n \log \exp\left(-\frac{1}{2\sigma^2} \|x - \mu\|^2\right) \\ &= \sum_{i=1}^n \left(-\frac{d}{2} \log(2\pi\sigma^2)\right) + \sum_{i=1}^n \left(-\frac{1}{2\sigma^2} \|x - \mu\|^2\right) \\ &= -\frac{nd}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^n \|x - \mu\|^2 = L \end{aligned}$$

Differentiate loglikelihood with respect to  $\mu$  and  $\sigma^2$  set to zero and solve for the respective parameters yields:

$$\begin{aligned} \hat{\mu} &= \frac{\sum_{i=1}^n x^{(i)}}{n} \\ \hat{\sigma}^2 &= \frac{\sum_{i=1}^n \|x^{(i)} - \mu\|^2}{nd} \end{aligned}$$

### 6.5 Gaussian Mixture Models

Is called SSoft Clustering"because it deals with probabilities not hard classification.

We have  $K$  clusters, each with own gaussian cloud  $N(x, \mu^{(j)}, \sigma_{(j)}^2)$ ,  $j = 1, \dots, K$ .  
Each Cluster gets own mixture-weight  $j \sim \text{Multinomial}(p_1, \dots, p_K)$   
Parameters of the mixture model are parameters of Multinomials and gaussians:

$$\theta = p_1, \dots, p_K; \mu^{(1)}, \dots, \mu^{(k)}, \sigma_{(j)}^2, \dots, \sigma_{(j)}^2)$$

Conditional probability of data-point given gaussian mixture:

$$P(x|\theta) = \sum_{i=1}^K p_j N(x, \mu^{(j)}, \sigma_{(j)}^2)$$

Conditional Likelihood of Training set given gaussian mixture:

$$P(S_n|\theta) = \prod_{j=1}^n \sum_{i=1}^K N(x, \mu^{(j)}, \sigma_{(j)}^2)$$

#### Observed Case:

We know to which mixture  $x$  belongs.

Indicator Variable is used to count the cases in which observation is part of a cluster  $\delta(j|i) = 1(x^{(i)} \text{ is assigned to } j)$ .

$$\begin{aligned} \sum_{i=1}^n [\sum_{j=1}^K \delta(j|i) \log(p_j N(x, \mu^{(j)}, \sigma_{(j)}^2))] &= \\ &= \sum_{j=1}^K [\sum_{i=1}^n \delta(j|i) \log(p_j N(x, \mu^{(j)}, \sigma_{(j)}^2))] \end{aligned}$$

Optimizing (according to MLE principle) yields:

$$\begin{aligned} \hat{n}_j &= \sum_{i=1}^n \delta(j|i) \\ \hat{p}_j &= \frac{\hat{n}_j}{n} \\ \hat{\mu}^{(j)} &= \frac{1}{\hat{n}} \sum_{i=1}^n \delta(j|i) \cdot x^{(i)} \\ \hat{\sigma}^2 &= \frac{1}{\hat{n}_j} \sum_{i=1}^n \delta(j|i) \|x^{(i)} - \mu^{(j)}\|^2 \end{aligned}$$

#### EM Algorithm (Unobserved Case):

We don't know to which mixture  $x$  belongs.

1. Randomly initialize  $\theta = p_1, \dots, p_K; \mu^{(1)}, \dots, \mu^{(k)}, \sigma_{(j)}^2, \dots, \sigma_{(j)}^2)$
2. E-Step:
  - (a) Calculate the softcount of a point (the probability of a cluster  $j$  given the point  $i$ :  $p(j|i) = \frac{p_j N(x^{(i)}; \mu^{(j)}, \sigma_{(j)}^2)}{p(x^{(i)}|\theta)}$ , where  $P(x|\theta) = \sum_{j=1}^K p_j N(x, \mu^{(j)}, \sigma_{(j)}^2)$
3. M-Step
  - (a) Use softcounts to calculate new parameters.

$$\begin{aligned} \hat{n}_j &= \sum_{i=1}^n p(j|i) \\ \hat{p}_j &= \frac{\hat{n}_j}{n} \\ \hat{\mu}^{(j)} &= \frac{1}{\hat{n}} \sum_{i=1}^n p(j|i) \cdot x^{(i)} \\ \hat{\sigma}_j^2 &= \frac{1}{\hat{n}_j} \sum_{i=1}^n p(j|i) (x^{(i)} - \mu^{(j)})^2 \end{aligned}$$

### 6.6 Reinforcement Learning

A Markov decision process (MDP) is defined by a set of states  $s \in S$  a set of actions  $a \in A$ ; Action dependent transition probabilities  $T(s, a, s') = P(s'|s, a)$ , so that for each state  $s$  and action  $a$ :  $\sum_{s' \in S} T(s, a, s') = 1$ .

Reward functions  $R(s, a, s')$  representing the reward for starting in state  $s$ , taking action  $a$  and ending up in state  $s'$  after one step. (The reward function may also depend only on  $s$ , or only  $s$  and  $a$ .) Therefore a Markov decision process is defined by  $MDP = \langle S, A, T, R \rangle$  MDPs satisfy the Markov property in that the transition probabilities and rewards depend only on the current state and action, and remain unchanged regardless of the history (i.e. past states and actions) that leads to the current state. Rewards collected after the  $n$ th step do not depend on the previous states  $s_1, s_2, \dots, s_{n-1}$   
**Markov properties:**

Rewards collected after the  $n$ th step do not depend on the previous actions  $a_1, a_2, \dots, a_n$   
**(Infinite horizon) discounted reward based utility**

$$\begin{aligned} U[s_0, s_1, \dots] &= R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) \dots = \\ &= \sum_{i=0}^{\infty} \gamma^i R(s_i) \text{ where } 0 \leq \gamma < 1 \\ &\leq \frac{R_{max}}{\gamma} \end{aligned}$$

#### Bellman Equations

$$V^*(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

$$Q^*(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

Q-value:  $Q(s, a)$  in state  $s$  take action  $a$  and act optimally afterwards. Policy  $\pi^*$ :  $s \rightarrow a$  is set of actions to maximize the expected reward for every state  $s$ .

$$\pi^*(s) = \arg\max_a (Q^*(s, a))$$

$$Q^*(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \max_{a'} Q(s', a')]$$

To find the policy two algos: Value iteration and Q-value iteration (look online).

### 6.7 Q value iteration by sampling

#### 6.8 A fully connected FFNN

Forward propagation

$$\begin{aligned} \text{Initialize the input layer } \mathbf{Z}^1 &= \mathbf{X}, \mathbf{A}^1 = \mathbf{Z}^1 \\ \text{Propagate all activity forward } \mathbf{Z}^l &= \mathbf{W}^l \mathbf{A}^{l-1} + \mathbf{B}^l, \mathbf{A}^l = f^l(\mathbf{Z}^l) \end{aligned} \quad (1)$$

$\mathbf{Z}^l = \mathbf{W}^l \mathbf{A}^{l-1} + \mathbf{B}^l$  is equivalent to

$$\begin{bmatrix} Z_1 \\ \vdots \\ Z_j \end{bmatrix} = \begin{bmatrix} W_{1,1} & \dots & W_{1,K} \\ \vdots & \ddots & \vdots \\ W_{j,1} & \dots & W_{j,K} \end{bmatrix} \begin{bmatrix} A_1 \\ \vdots \\ A_K \end{bmatrix} + \begin{bmatrix} B_1 \\ \vdots \\ B_j \end{bmatrix} \quad (2)$$

Where layer  $l$  has  $J$  nodes and layer  $l-1$  has  $K$  nodes.  
Backward propagation

Assuming that the loss function is  $\mathcal{L} = 0.5(\mathbf{Y} - \mathbf{A}^L)^2$   
Where  $\mathbf{Y}$  is the known output corresponding to the input  $\mathbf{X}$

Calculate the final error  $\nabla_a \mathcal{L} = \mathbf{A}^L - \mathbf{Y}$

Initialize the back propagation  $\delta^L = \nabla_a \mathcal{L} * f^{L'}(\mathbf{Z}^L)$

Backpropagate the error  $\delta^l = \left[ \mathbf{W}^{l+1 T} \delta^{l+1} \right] * f^{l'}(\mathbf{Z}^l), l \in [2, L-1]$

Calculate the gradient of weights  $\frac{\partial \mathcal{L}}{\partial \mathbf{W}^l} = \delta^l \mathbf{A}^{l-1 T}, l \in [2, L]$

Equivalently  $\frac{\partial \mathcal{L}}{\partial W_{j,k}^l} = a_k^{l-1} \delta_j^l, l \in [2, L]$

Calculate the gradient of bias weights  $\frac{\partial \mathcal{L}}{\partial \mathbf{B}^l} = \delta^l, l \in [2, L]$

Equivalently  $\frac{\partial \mathcal{L}}{\partial B_j^l} = \delta_j^l, l \in [2, L]$

Update ( $\eta$  is a hyperparameter that is not learned by the FFNN)

Update weights  $\mathbf{W}^l = \mathbf{W}^l - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{W}^l}, l \in [2, L]$

Update bias weights  $\mathbf{B}^l = \mathbf{B}^l - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{B}^l}, l \in [2, L]$

### 6.9 ATE abd Selection Bias

When we observe a group of subjects, some of whom have been treated and some not, we calculate

$$\begin{aligned} E(Y_i(1)|W_i = 1) - E(Y_i(0)|W_i = 0) &\iff \\ [E(Y_i(1)|W_i = 1) - E(Y_i(0)|W_i = 1)] + \\ E(Y_i(0)|W_i = 1) - E(Y_i(0)|W_i = 0)] &\iff \end{aligned} \quad (5)$$

Average Treatment Effect + Selection Bias

### 6.10 Fixed effects model

G(factor(admin)) creates a dummy variable, one per region, and includes the dummy variable in the regression

```
library("lme4")
model2 <- felm(sex ~ teasing + post + teapost +
G(factor(admin)), data = qiandata)
summary(model2)
```

### 6.11 Instrument variable regression

To use the ivreg() function directly one needs to setup the call as

```
ivreg (y ~ exogenous_vars + endogenous_var |
exogenous_vars + instrument_var)
```

### 6.12 Fisher exact test in R

```
library(permm)
rm(list = ls())
cough_severity = matrix(c(3, 5, 0, 4, 0, 1), nrow=6,
byrow=TRUE, ncol=1, byrow=TRUE)
n_treatment = 3
n_control = 3
n = n_treatment + n_control

permutations = chooseMatrix(n, n_treatment)
num_permutations = nrow(permutations)

treatment_avg = (1/n_treatment)* permutations %*%
by_cough_severity

control_avg = (1/n_control)*(1-permutations)%*%
by_cough_severity

test_statistic = abs(treatment_avg-control_avg)

actual_obs_row = 20
observed_test_statistic = test_statistic[actual_obs_row]

by_chance = (test_statistic >= observed_test_statistic)
total_by_chance = sum(by_chance)
```