# Naive Bayes Classifier

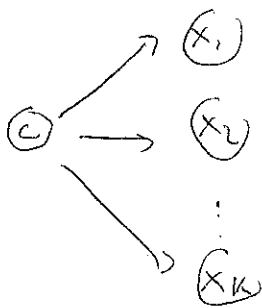<u>D</u>o <u>this</u> <u>earlier</u>

Have $\wedge$ observation X which we assign to $1$ of C classes : $1 \cdots C$.

(k-tuple)

X can be real-valued, binary, ...

Naive Bayes (NB) models

$$q(x,c) = q(c) \, q(x|c) \overset{NB}{=} q(c) \prod_{K=1}^{K} q(x_k|c)$$



$\Longleftrightarrow$    $X_1 \cdots X_k$ cond. indep. given $c$.

<u>Note</u>  If we don't assume cond. indep. must learn k-dimen. dist for each class.  C.I means we learn k 1-d dists. for each class.

<u>Ex</u>  Suppose  $x | c \sim N(\mu_c, \Sigma_c)$   $\frac{k^2}{?} + k$ parms for each class

vs.    $x_k | c \sim N(\mu_{kc}, \sigma_{kc}^2)$   $2k$ parms   " "   " .

$X_1 \cdots X_k$ cond. indep. $| c$

Fewer parameters means more accurate estimation, ~~XXXXXX~~

Restrictive assump. (if wrong) means less accurate model

This is fundamental tradeoff of ML.

Discrete

~~X~~ Features + Naive Bayes

Say $X = \begin{pmatrix} x_1 \\ \vdots \\ x_k \end{pmatrix}$ has ~~X~~ $x_k \in \{1, 2, \cdots M\}$

Need (according to NB) $p(x_k = m \mid c)$ 

$k = 1 \cdots K$
$m = 1 \cdots M$
$c = 1 \cdots C$.

Let $X$ be data matrix
$n \times k$

1st feat    $k$th feat
  $\downarrow$       $\downarrow$

$X = \begin{pmatrix} x_{11} & \cdots & x_{1k} \\ x_{21} & \cdots & x_{2k} \\ & \vdots & \\ x_{n1} & \cdots & x_{nk} \end{pmatrix}$  $\leftarrow$ 1st obs.

$\leftarrow$ $n$th obs.

$C = \begin{pmatrix} c_1 \\ \vdots \\ c_n \end{pmatrix}$

How to estimate $p(x_k = m \mid c)$?

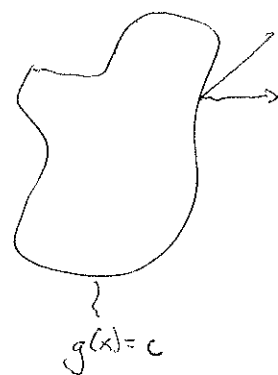Let $n_{kmc} = \left| \{ i : x_{ik} = m, c_i = c \} \right|$

( Estimating Probabilities for Discrete Events. )

we expt. that has $M$ outcomes $\{1 \cdots M\}$ with probs $q_1 \cdots q_M$ $\left(\sum q_n = 1\right)$
Want to estimate $q_1 \cdots q_M$.

Have $n$ independent observations of expt. $X_1 \cdots X_n$    $X_i \in \{1 \cdots M\}$.
The maximum likelihood est. satisfies ⌐ Let $n_m = |\{i : x_i = m\}|$ $\sum_{m=1}^{M} n_m = n$

$$\hat{q}_1 \cdots \hat{q}_M = \arg\max_{\substack{q_1 \cdots q_M \\ s.t. \ \sum q_n = 1}} \prod_{i=1}^{n} q_{x_i} = \arg\max_{\substack{q_1 \cdots q_M \\ s.t. \ \sum q_n = 1}} \prod_{m=1}^{M} q_m^{n_m}$$

$$= \arg\max_{\substack{q_1 \cdots q_M \\ s.t. \ \sum q_n = 1}} \sum_{m=1}^{M} n_m \log q_m$$

Lagrange multipliers

$$\nabla \sum_{m=1}^{M} n_m \log q_m = \lambda \nabla \sum q_n = 1$$

$$\Longleftrightarrow$$

$$\begin{pmatrix} \frac{n_1}{q_1} \\ \vdots \\ \frac{n_M}{q_M} \end{pmatrix} = \lambda \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \implies$$

① $q_m = n_m / \lambda$

② $q_m = n_m / n$

$g(x) = c$

( Answer is obvious guess but MLE supports guess )

# Discrete Features + Naive Bayes

Have classification problem with data $\underset{n \times k}{X}$ and class $\underset{n \times 1}{c}$

$$X = \begin{pmatrix} X_{11} & X_{12} & \cdots & X_{1k} \\ X_{21} & X_{22} & \cdots & X_{2k} \\ \vdots & & & \\ X_{n1} & X_{n2} & \cdots & X_{nk} \end{pmatrix} \begin{matrix} \leftarrow \text{obs } 1 \\ \\ \\ \leftarrow \text{obs } n \end{matrix} \qquad C = \begin{pmatrix} C_1 \\ \vdots \\ C_n \end{pmatrix}$$

$\uparrow$ feat $\underline{1}$ $\qquad \uparrow$ Feat $k$

$c_i \in \{1 \cdots C\} \qquad X_{ik} \in \{1, \cdots M\}$

NB requires we have $\hat{p}(x_k = m \mid c)$ $\qquad \begin{matrix} k=1 \cdots K \\ m=1 \cdots M \\ c=1 \cdots C \end{matrix}$

Let $n_{kmc} = |\{i : X_{ik} = m, c_i = c\}| = $ # times get class $c$ with $k^{th}$ feat $= m$.

By prev. argument

$$\hat{p}(x_k = m \mid c) = \frac{n_{kmc}}{\sum_{m'=1}^{M} n_{km'c}}$$

And also

$$\hat{p}(c) = \frac{|\{i : c_i = c\}|}{n}$$

<u>view</u>

$$\sigma(t) = \frac{1}{1+e^{-t}} \qquad \text{Model} \quad q(c=1|x) = \sigma(w^t x)$$

Have data vectors $x_1, x_2, \dots x_n$     We know

$$\frac{\partial \log q(c_1 \cdots c_n | x_1 \cdots x_n)}{\partial w_j} = \sum_{i=1}^{n} (c_i - \sigma(w^t x_i)) x_{ij} \iff \frac{\nabla \log q(c_1 \cdots c_n | x_1 \cdots x_n)}{\phantom{x}}$$
$$= \sum (c_i - \sigma(w^t x_i)) x_i$$

Writing
$$\underset{n \times k}{X} = \begin{pmatrix} x_{11} & \cdots & x_{1k} \\ x_{21} & \cdots & x_{2k} \\ \vdots & & \\ x_{n1} & \cdots & x_{nk} \end{pmatrix} \qquad c = \begin{pmatrix} c_1 \\ \vdots \\ c_n \end{pmatrix} \qquad q = \begin{pmatrix} \sigma(w^t x_1) \\ \vdots \\ \sigma(w^t x_n) \end{pmatrix}$$

$$\nabla \log q(c_1 \cdots c_n | x_1 \cdots x_n) = X^t (c - q)$$

$$\frac{\partial^2 \log q(c_1 \cdots c_n | x_1 \cdots x_n)}{\partial w_j \partial w_\ell} = -\sum_{i=1}^{n} \sigma(w^t x_i)(1 - \sigma(w^t x_i)) x_{ij} x_{i\ell}$$

$$= -(X^t C X)_{j\ell}$$

$$C = \begin{pmatrix} \sigma(w^t x_1)(1-\sigma(w^t x_1)) & & O \\ & \ddots & \\ O & & \sigma(w^t x_n)(1-\sigma(w^t x_n)) \end{pmatrix} \qquad \text{ie.} \quad -X^t C X \quad \text{is Hessian}$$

NR:   $w^{new} = w^{old} + H^{-1} \nabla f(w^{old}) = w^{old} + (X^t C X)^{-1} X^t (c - q)$

[ NB. $C$ and $q$ depend on $w^{old}$ so must be recomputed each iter

# Regression

In classification data are labeled samples:

$$(x_1, c_1), (x_2, c_2) \dots (x_n, c_n)$$

where $\{x_i\}$ are vectors usually $\in \mathbb{R}^d$ and $c_i \in \underbrace{\{1, \dots c\}}_{\text{our classes}}$

In regression have $(x_1, y_1) \dots (x_n, y_n)$ where

$$x_i \in \mathbb{R}^d \ (\text{as before}) \text{ but } y_i \in \mathbb{R} \ (\text{or } \mathbb{R}^k)$$

The $\{x_i\}$ are the predictors and the $\{y_i\}$ are the response.

## Formulation with Loss Function

Let $\hat{y} = \hat{y}(x)$ be prediction of $y$ based on $x$

Suppose we adopt loss function $L(y, \hat{y}) = (y - \hat{y})^2$

giving "cost" for estimating $\hat{y}$ when truth is $y$.

Observe $x$ and seek $\hat{y}$ that minimizes expected loss

$$E L = \int (y - \hat{y})^2 \, p(y|x) \, dy$$

Know $\hat{y} = E(y|x) = \int y \, p(y|x) \, dy$.

so $\hat{y}(x) = E(y|x)$ is obvious choice for regression.

$\hat{y}(x) = E(y|x)$ is known as regression function

# Linear Regression

In linear regression predict $y$ as linear fn of $x$

$$\hat{y} = \hat{y}(x) = w^t x = w_1 x_1 + w_2 x_2 + \cdots w_k x_k .$$

## Ex

$$x = \underset{x_1}{\text{amount of education}}, \underset{x_2}{\text{salary of parents}}, \underset{x_3}{\text{age}}$$

$y =$ salary of individual

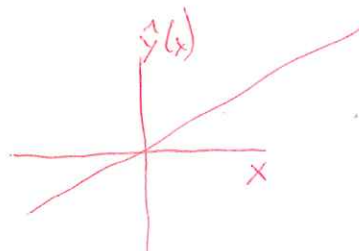Estimate $\hat{y} = w_1 x_1 + w_2 x_2 + w_3 x_3 = w^t x$ .

In linear regression common to augment observed predictors with other variables derived from observations.
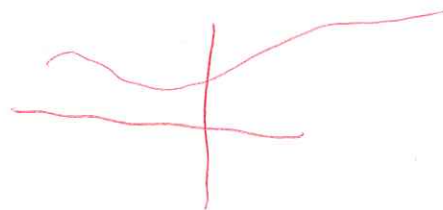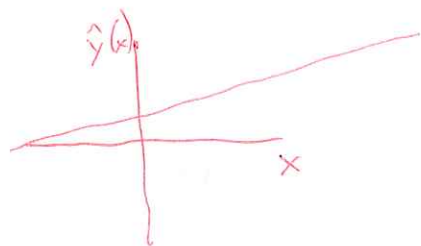
## Ex

Suppose data look like



Straight linear regression requires $\hat{y}(x) = wx$

(line through origin)

But could view predictors as $1, x$ so $\hat{y}(x) = w_0 1 + w_1 x$



or perhaps $1, x, x^2, x^3 \cdots x^k$ so $\hat{y}(x) = w_0 1 + w_1 x + w_2 x^2 + \cdots + w_k x^k$

This is still linear regression !

# Geometric View of Regression

Have data $(x_1, y_1) \cdots (x_n, y_n)$    $x_i \in \mathbb{R}^k$, $y_i \in \mathbb{R}$.

$X_i$ includes whatever features we derive from obs including $\underline{1}$.

Use linear prediction:   $\hat{y}_i = w^t x_i$   and want to minimize sum of squared errors (SSE) between $\hat{y}_i$ and $y_i$

$$\hat{w} = \arg\min_{w} \sum_{i=1}^{n} (y_i - w^t x_i)^2$$

## Another View

$$\text{Let} \quad X = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1k} \\ x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & & & \\ x_{n1} & x_{n2} & \cdots & x_{nk} \end{pmatrix} \qquad y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \qquad X \text{ is } \underline{\text{data matrix}}.$$
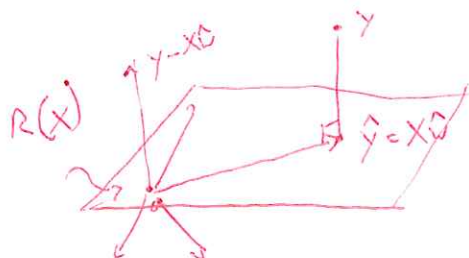
(1st obs → top row, nth obs → bottom row; 1st feat, $n^{th}$ feat columns)

$$\hat{w} = \arg\min_{w} \sum_{i=1}^{n} \left(y_i - (Xw)_i\right)^2 = \arg\min_{w} \| y - Xw \|^2$$

## Geometric Picture

Write $R(X) = \{ Xw : w \in \mathbb{R}^k \}$ = range of X.

Seek $\hat{y} \in R(X)$ s.t. $\| y - \hat{y} \|^2$ minimized.



Pictorially $y - \hat{y}$ should be orthogonal to $R(X)$

$$(Xw, y - X\hat{w}) = 0 \quad \forall w$$

$$\Longleftrightarrow$$

$$(w, X^t(y - X\hat{w})) = 0 \quad \forall w \quad \Longleftrightarrow \quad X^t(y - X\hat{w}) = 0$$

(right margin, bottom to top:)

$$X^t X \hat{w} = X^t y \quad (\text{normal eqns})$$

$$\Uparrow$$

$$\hat{w} = (X^t X)^{-1} X^t y.$$

Easy to remember version

Would like to solve $\quad Xw = y$
$\qquad\qquad\qquad n\times k \quad k\times 1 \quad n\times 1$

Think of

$$\begin{pmatrix} X \end{pmatrix} \begin{pmatrix} w \end{pmatrix}_{k\times 1} = \begin{pmatrix} y \end{pmatrix}$$

$n\times k \qquad\qquad\qquad n\times 1$

Could have many (n)
='ns, but few (k)
unknowns so can't
solve this usually.

But, can solve $\quad X^t X w = X^t y \quad$ (Normal ='ns)

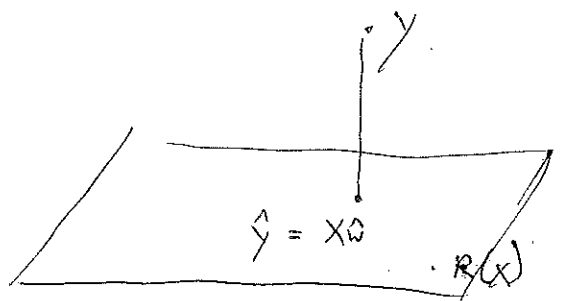Ex   Simple-linear-regression.r

Projection

We saw

$\hat{w} = \arg\min_w \| y - Xw \|$

given by

① $\quad X^t X \hat{w} = X^t y \quad$ (Normal ='ns)

② $\quad \hat{w} = (X^t X)^{-1} X^t y$



$\hat{y} = X\hat{w}$ is "projection" of
$y$ onto $R(X)$.

Projection Matrices

If  L  is a linear space (eg plane or line containing origin)
the "projection" of  y  onto $L^{\wedge}$ is $\overset{P_L y}{}$ closest point in L to y.

A projection matrix  P  has

① $P^2 = P$  (idempotent)

② $P^t = P$  (symmetric)

Easy to see  $X(X^tX)^{-1}X^t$  satisfies ① + ②

① $\left(X(X^tX)^{-1}X^t\right)\left(X(X^tX)^{-1}X^t\right) = X(X^tX)^{-1}X^t$

② $\left(X(X^tX)^{-1}X^t\right)^t = X(X^tX)^{-1^t}X^t = X(X^tX)^{-1}X^t$

Since $R\left(X(X^tX)^{-1}X^t\right) = R(X)$      $X(X^tX)^{-1}X^t = P_{R(X)}$

.Thm of Pythagoras for proj

Let  P  be projection. For any  y

$y = Py + (I-P)y \implies$

$\|y\|^2 = \|Py + (I-P)y\|^2 = \left(Py + (I-P)y, Py + (I-P)y\right)$

$= \|Py\|^2 + \|(I-P)y\|^2$

Note $(Py, (I-P)y) = (y, P(I-P)y) = (y, 0y) = 0$

$\|y\|^2 = \|\hat{y}\|^2 + \|y-\hat{y}\|^2$

For our special case of $\hat{y} = X\hat{b}$

$\langle \hat{y} = X\hat{b}$

# Statistical View of Regression

( ~~Fixed~~ Data $(x_1, y_1) \dots (x_n, y_n)$ $x_i \in \mathbb{R}^k$, $y_i \in \mathbb{R}$

### Model

$$y_i = w^t x_i + \varepsilon_i \qquad i = 1 \dots n$$

where $w$ is vector of unknowns $w \in \mathbb{R}^k$ $\quad \left[ w, \sigma^2 \text{ parameters} \right]$

$\qquad \varepsilon_1 \dots \varepsilon_n \overset{iid}{\sim} N(0, \sigma^2)$ $\qquad \sigma^2$ unknown.

### Equivalently

$$X_{n \times k} = \begin{pmatrix} x_{11} & \cdots & x_{1k} \\ x_{21} & \cdots & x_{2k} \\ \vdots & & \vdots \\ x_{n1} & \cdots & x_{nk} \end{pmatrix} \qquad y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \qquad w = \begin{pmatrix} w_1 \\ \vdots \\ w_k \end{pmatrix} \qquad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

$$y = Xw + \varepsilon \qquad\qquad \varepsilon \sim N(0, \sigma^2 I)$$

### Defn

An estimator, $\hat{\Theta}$, for parm $\Theta$ is unbiased if $E\hat{\Theta} = \Theta$

(on average estimate is correct)

### Ex

If $X_1 \dots X_n$ are sequence of indep. S-F trials

$$P(X_i = 1) = p$$
$$P(X_i = 0) = 1 - p$$

We saw the MLE for $p$ was $\hat{p} = \frac{1}{n} \sum X_i$.

Easy to see $\hat{p}$ is unbiased $E\hat{p} = ?$

$E X_i = 1 \cdot p + 0 \cdot (1-p) = p$

$E\hat{p} = E \frac{1}{n} \sum X_i = \frac{1}{n} \sum E X_i = \frac{1}{n} \sum p = \frac{1}{n} \cdot np = p$

We have used $\hat{w} = (X^t X)^{-1} X^t y$ as estimate for $w$.

In fact $\hat{w}$ is unbiased for $w$.

### Have

$$y = Xw + \varepsilon \qquad \varepsilon \sim N(0, \sigma^2 I)$$

$$E\hat{w} = E(X^t X)^{-1} X^t y = E(X^t X)^{-1} X^t (Xw + \varepsilon)$$

$$= E(X^t X)^{-1} X^t X w + E(X^t X)^{-1} X^t \varepsilon$$

$$= w + (X^t X)^{-1} X^t \underbrace{E \varepsilon}_{0}$$

$$= w \qquad \Longrightarrow \qquad \hat{w} \text{ unbiased for } w.$$

### Variance of Unbiased Estimator

If $\hat{\theta}$ is unbiased for $\theta$, then

$$V(\hat{\theta}) = E(\hat{\theta} - \theta)^2 = \text{expected sq. error}.$$

If $\hat{\theta}_1$ and $\hat{\theta}_2$ are UE's for $\theta$ and $V(\theta_1) < V(\theta_2)$ then $\hat{\theta}_1$ has less sq error on average, and is better in this regard.

This variance often used as goodness meas. for UEs.

### Gauss-Markov Thm

Informally, $\hat{w} = (X^t X)^{-1} X^t y$ has smallest variance (sq. error) of all UE's.

More precisely, let $\tilde{w}$ be a U.E. for $w$, and $\alpha \in \mathbb{R}^k$. Then

$$V(\alpha^t \hat{w}) \le V(\alpha^t \tilde{w})$$

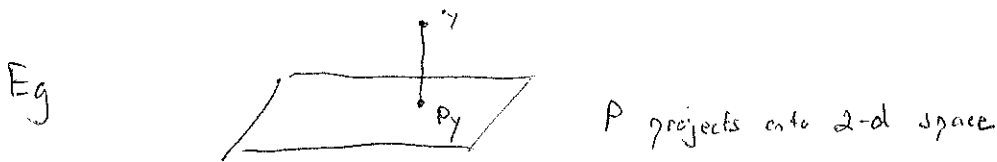For ex, if $\alpha^t = (0 \dots \overset{i^{th}}{1} \dots 0)$ this says $V(\hat{w}_i) \le V(\tilde{w}_i)$

# Estimating $\sigma^2$

① If $x \sim N(0, \sigma^2)$   $E x^2 = \int x^2 N(x; 0, \sigma^2) dx = \sigma^2$

If $\underset{n \times 1}{\varepsilon} \sim N(0, \underset{n \times n}{\sigma^2 I})$

$E \| \varepsilon^2 \|^2 = E \left( \varepsilon_1^2 + \varepsilon_2^2 + \cdots + \varepsilon_n^2 \right) = n \sigma^2$

② Fact   Suppose $\underset{n \times n}{P}$ is projection matrix proj onto $k$-dim space

Eg



$P$ projects onto 2-d space

Then $E \| P \varepsilon \|^2 = k \sigma^2$

Since   $\| \varepsilon \|^2 = \| P \varepsilon \|^2 + \| (I - P) \varepsilon \|^2$

Taking expectations gives:

$n \sigma^2 = k \sigma^2 + E \| (I - P) \varepsilon \|^2$

$\implies E \| (I - P) \varepsilon \|^2 = (n - k) \sigma^2$

③ In regression model $y = Xw + \varepsilon$   let $P_{R(x)}$ be proj onto $R(x)$

$E \| y - \hat{y} \|^2 = E \| (I - P_{R(x)}) y \|^2 = E \| (I - P_{R(x)}) Xw + \varepsilon \|^2$

$\qquad = E \| (I - P_{R(x)}) \varepsilon \|^2 = (n - k) \sigma^2$

$\implies E \frac{\| y - \hat{y} \|^2}{n - k} = \sigma^2$

We will use   $\hat{\sigma}^2 = \frac{\| y - \hat{y} \|^2}{n - k}$   as our (unbiased) estimate for $\sigma^2$.

Do regression_variance.r

# Overfitting

Have data $(x_1, y_1) \ldots (x_n, y_n)$  $x_i \in \mathbb{R}^{*d}$, $y_i \in \mathbb{R}$

**Ex**  Want to predict price of Apple stock on particular day.

Choose relevant predictors

① Overall consumer spending

② Advertising expenditures of Apple

③ Investor confidence in tech sector

④ Price of labor is China

$$X = \begin{array}{c} \text{day 1} \\ \vdots \\ \text{day } n \end{array} \overset{\text{var 1} \quad \cdots \quad \text{var } d}{\left( \phantom{\begin{matrix} a \\ a \\ a \\ a \\ a \end{matrix}} \right)}$$

Measure variables and get  $SSE = \| y - \hat{y} \|^2 = \| y - X\hat{w} \|^2$

SSE doesn't seem small enough ( predictions, $\hat{y}$, not close to $y$ )
so add new predictors

ⓐ Rainfall in Ecuador on each day

ⓑ Price of Barley futures

ⓒ Dist between closest pair of Jupiter's moons

⋮   ⋮

$$X = \begin{array}{c} \text{day 1} \\ \\ \\ \text{day } n \end{array} \overset{\text{Vars} \atop 1 \; 2 \cdots d, a, b, c, \cdots}{\left( \phantom{\begin{matrix} a \\ a \\ a \\ a \\ a \end{matrix}} \right)}$$

Since new predictors are independent of Apple's stock price, so shouldn't help.
But SSE continually decreases as irrelevant predictors added.
This <u>seems</u> like good news, but with <u>new data</u>  $(x_{n+1}, y_{n+1}) \ldots (x_{2n}, y_{2n})$
learned model predicts poorly. That is

$X_T, y_T$ are original data, $\hat{w}_T$ learned weights  $\hat{w}_T = (X_T^t X)^{-1} X_T^t y_T$
$X_N, y_N$ new data  $SSE_N = \| y_N - X_N \hat{w}_T \|^2$ is high. This is known as overfitting

# Ways to Avoid Overfitting

① Variable selection: Use only subset of vbles.

Use notation

$$X = \begin{pmatrix} x_{11} & \cdots & x_{1d} \\ x_{21} & \cdots & x_{2d} \\ \vdots & & \vdots \\ x_{n1} & \cdots & x_{nd} \end{pmatrix} \qquad Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

$$X(j) = j^{th} \text{ col of } X \qquad \begin{pmatrix} x_{1j} \\ x_{2j} \\ \vdots \\ x_{nj} \end{pmatrix}$$

$$X(j_1, j_2) = \text{cols } j_1, j_2 \qquad \begin{pmatrix} x_{1j_1} & x_{1j_2} \\ \vdots & \vdots \\ x_{nj_1} & x_{nj_2} \end{pmatrix}_{n \times 2}$$

$$X(j_1, j_2, j_3) = \text{cols } j_1, j_2, j_3 \text{, etc.}$$

Let $j_1$ be best single predictor.

$$j_1 = \arg\min_j \| y - X(j) \cancel{\phantom{xxx}} (X^t(j) X(j))^{-1} X^t(j) y \|^2$$

Choose $2^d$ predictor, $j_2$ to be best vble in addition to $j_1$

$$j_2 = \arg\min_j \| y - X(j_1, j) (X^t(j_1, j) X(j_1, j))^{-1} X^t(j_1, j) y \|^2$$

$$j_3 = \arg\min_j \| y - X(j_1, j_2, j) (X^t(j_1, j_2, j) X(j_1, j_2, j))^{-1} X^t(j_1, j_2, j) y \|^2 \cdots$$

Can show that under non-pathological conds.

$SSE_0 > SSE_1 > SSE_2 \cdots$

How far should we go?

Add vbles as long as

$SSE_{k-1} - SSE_k > \text{threshold}.$

Let

$$SSE_0 = \|y\|^2$$

$$SSE_1 = \| y - X(j_1)(X^t(j_1) X(j_1))^{-1} X^t(j_1) y \|^2$$

$$SSE_2 = \| y - X(j_1, j_2)(X^t(j_1, j_2) X(j_1, j_2))^{-1} X^t(j_1, j_2) y \|^2$$

$\cdots$

# Ridge Regression

( Initial formulation of regression: $\hat{w} = \arg\min_w \|y - Xw\|^2$.

Have seen this formulation is prone to overfitting with many predictors.

Ridge regression penalizes complex fits of the data:

$$\hat{w}_{Ridge} = \arg\min_w \underbrace{\|y - Xw\|^2}_{\text{data fit}} + \underbrace{\lambda\|w\|^2}_{\text{complexity penalty}}. \qquad \lambda > 0$$

Can show $\hat{w}_{Ridge} = \left(X^t X + \lambda I\right)^{-1} X^t y$

When $\lambda = 0$ get old solution $\hat{w}$. As $\lambda$ increases $\hat{w}_{Ridge}$ "shrinks" to $0$.

( Note Ridge Regression useful when $\underbrace{X^t X}_{k \times k}$ is singular (eg. $k > n$)
$\underset{k \times n}{X^t}\underset{n \times k}{X}$

since $\left(X^t X + \lambda I\right)$ always invertible.

<u>How to Choose $\lambda$</u>? ( Cross validation )

Suppose we divide data into "training set" and "validation set"

$$X = \left(\begin{array}{ccc} x_{11} & x_{12} & \cdots & x_{1k} \\ x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & & & \\ \hline \vdots & & & \\ x_{n1} & x_{n2} & \cdots & x_{nk} \end{array}\right) \begin{array}{l} X_T \\ \\ X_V \end{array} \qquad y = \left(\begin{array}{c} y_1 \\ y_2 \\ \vdots \\ \hline \\ y_n \end{array}\right) \begin{array}{l} y_T \\ \\ y_V \end{array} = \left(\begin{array}{c} y_T \\ y_V \end{array}\right)$$

$$= \left(\begin{array}{c} X_T \\ X_V \end{array}\right)$$

$$\hat{w}_T(\lambda) = \arg\min_w \|y - X_T w\|^2 + \lambda\|w\|^2$$

$$= \left(X_T^t X_T + \lambda I\right)^{-1} X_T^t y$$

Choose $\lambda$ giving best performance on validation set.

$$\hat\lambda = \arg\min_\lambda \|y_V - X_V \hat{w}_T(\lambda)\|$$

# LASSO Regression

$$\hat{w}_{LASSO} = \arg\min_{w} \|y - Xw\|^2 + \lambda \|w\|$$

no square.

① As before $\lambda = 0$ gives $\hat{w} = (X^t X)^{-1} X^t y$.

② As $\lambda$ increases some components of $\hat{w}_{LASSO}$ driven to $0$.
   thus LASSO acts like variable selection while $\lambda$ controls # of vbles.

③ Computation of $\hat{w}_{LASSO}$ more complicated + won't discuss.

# Gaussian Mixtures + EM Algorithm

## Mixtures

① Choose randomly from $K$ classes $\{1 \cdots k\}$ with probs $\Pi_1 \cdots \Pi_k$.

② If choose $k^{th}$ class sample $X$ from $X \sim q_k(x)$

$$q(x) = \sum_{k=1}^{k} \Pi_k q_k(x)$$

For Gaussian Mixture Model (GMM) class cond. dists $(q_k(x))$ are multivariate normal

$$q(x) = \sum_{k=1}^{k} \Pi_k N(x; \mu_k, \Sigma_k) = \sum_{k=1}^{K} \Pi_k \frac{1}{(2\pi)^{d/2} |\Sigma_k|^{1/2}} e^{-\frac{1}{2}(x-\mu_k)\Sigma_k^{-1}(x-\mu_k)}$$

Suppose have data :

$$\underset{n \times d}{X} = \begin{pmatrix} x_{11} & \cdots & x_{1d} \\ x_{21} & \cdots & x_{2d} \\ \vdots & & \vdots \\ x_{n1} & \cdots & x_{nd} \end{pmatrix}$$

Want to estimate GMM $\{\Pi_k, \mu_k, \Sigma_k\}_{k=1}^{K}$

## Intuition

Let $c_i \in \{1 \cdots k\}$ be class $x_i$ comes from. (unknown)

If the $\{c_i\}$ known then let

$$n_k = |\{i : c_i = k\}| \quad ; \quad \hat{\mu}_k = \frac{1}{n_k} \sum_{i=1}^{n} x_i \mathbb{1}_{c_i = k} \quad ;$$

$$\hat{\Sigma}_k = \frac{1}{n_k} \sum_{i=1}^{n} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^t \mathbb{1}_{c_i = k}$$

where $\mathbb{1}_{c_i = k} = \begin{cases} 1 & \text{if } c_i = k \\ 0 & \text{o.w.} \end{cases}$

But we don't know $\{c_i\}$. Do know

$$p(c_i = k \mid x_i) \overset{Bayes}{=} \frac{\pi_k \, \mathcal{N}(x_i; \mu_k, \Sigma_k)}{\sum\limits_{k'=1}^{K} \pi_{k'} \, \mathcal{N}(x_i; \mu_{k'}, \Sigma_{k'})} \doteq \gamma_{ik}$$

$\gamma_{ik}$ is "responsibility" of $k^{th}$ class for $i$th sample

$$\sum_{k=1}^{K} \gamma_{ik} = 1 \qquad n = \sum_{i=1}^{n} \sum_{k=1}^{K} \gamma_{ik} = \sum_{k=1}^{K} \overbrace{\left( \sum_{i=1}^{n} \gamma_{ik} \right)}^{\;\doteq\, n_k}$$

$$\sum_{i=1}^{n} \cancel{\sum \gamma_{ik} \, n_k} \qquad\qquad = \sum_{k=1}^{K} \cancel{\gamma_{ik}} \, n_k$$

$n_{k}$ is # $\overset{samples}{\wedge}$ attributed to $k^{th}$ class (not an integer)

Idea: treat $x_i$ as $\gamma_{i1} \overset{of\ a}{\wedge}$ samples from class $1$

$\gamma_{i2}$ " " " $2$

$\vdots$

$\gamma_{ik}$ " " " $k$

$$\hat{\pi}_k = \cancel{} \; \frac{n_k}{n}$$

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{i=1}^{n} \gamma_{ik} \, x_i$$

$$\hat{\Sigma}_k = \frac{1}{n_k} \sum_{i=1}^{n} \gamma_{ik} (x_i - \mu_k)(x_i - \mu_k)^t$$

### Algorithm

① Init $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^{K}$

② Compute $\gamma_{ik} = p(c_i = k \mid x_i)$
and $n_k = \sum_{i=1}^{n} \gamma_{ik}$

③ (Re)Estimate $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^{K}$
from updates.

④ Go to ② until convergence.

Look at gmm.em.r

# Revisiting GMM Algorithm Viewed as Max Likelihood

$$X = \begin{pmatrix} x_{11} & \cdots & x_{1d} \\ x_{21} & \cdots & x_{2d} \\ \vdots & & \vdots \\ x_{n1} & \cdots & x_{nd} \end{pmatrix}$$

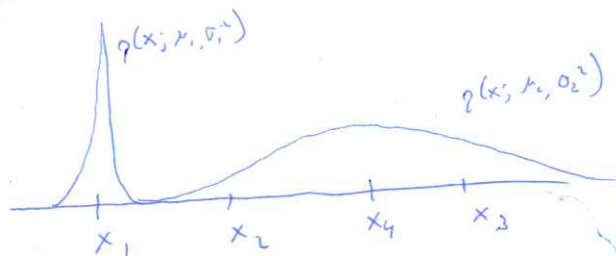Want to fit GMM to $X$ with $K$ - mixture components.

Have

$$p(x; \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^{k}) = \prod_{i=1}^{n} \sum_{k=1}^{K} \pi_k N(x_i; \mu_k, \Sigma_k)$$

Want to estimate params $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^{k}$ by MLE.

## Aside

Problem not well-posed since degenerate solutions give arbitrarily-high likelihood.

Consider. 1-d case with $K = 2$



If we take $\mu_1 = x$ and let $\sigma_1^2 \downarrow 0$ then

$$p(x_1; \mu_1, \sigma_1^2) = \frac{1}{(2\pi\sigma_1^2)^{1/2}} e^{-\frac{1}{2\sigma_1^2}(x_1 - \mu_1)^2} = \frac{1}{(2\pi\sigma_1^2)^{1/2}}$$

so $p(x_1; \mu_1, \sigma_1^2) \xrightarrow[\sigma_1^2 \downarrow 0]{} \infty$

Thus we seek non-degenerate solutions.

Now

$$\log p\left(X; \{\Pi_k, \mu_k, \Sigma_k\}_{k=1}^{k}\right) = \sum_{i=1}^{n} \log \sum_{k=1}^{K} \Pi_k \, \mathcal{N}(x_i; \mu_k, \Sigma_k)$$

could optimise by differentiating w.r.t. $\mu_k$ and setting to $0$

$$0 = \nabla_{\mu_k} \log p\left(X; \{\Pi_{k'}, \mu_{k'}, \Sigma_{k'}\}_{k'=1}^{K}\right)$$

$$= \sum_{i=1}^{n} \frac{\Pi_k \, \mathcal{N}(x_i; \mu_k, \Sigma_k) \, \Sigma_k^{-1} \, \cancel{(\text{///})} \, (\mu_k - x_i)}{\sum_{k'=1}^{K} \Pi_{k'} \, \mathcal{N}(x_i; \mu_{k'}, \Sigma_{k'})}$$

$$= \sum_{i=1}^{n} \gamma_{ik} \, \Sigma_k^{-1} (\mu_k - x_i)$$

$$= \Sigma_k^{-1} \sum_{i=1}^{n} \gamma_{ik} (\mu_k - x_i)$$

$$\implies \sum_{i=1}^{n} \gamma_{ik} (\mu_k - x_i) = 0$$

$$\implies \sum_{i=1}^{n} \gamma_{ik} \mu_k = \sum_{i=1}^{n} \gamma_{ik} \mu_k = n_k \mu_k$$

$$\implies \hat{\mu}_k = \frac{1}{n_k} \sum_{i=1}^{n} \gamma_{ik} x_i \qquad (\text{as before})$$

Can also show, differentiating w.r.t $\Sigma_k$ we get

$$\hat{\Sigma}_k = \frac{1}{n_k} \sum_{i=1}^{n} \gamma_{ik} (x_i - \hat{\mu}_k)(x_i - \mu_k)^t \qquad (\text{as before})$$

To estimate $\{\Pi_k\}$ since $\sum \Pi_k = 1$ need Lagrange:

$$\nabla_\Pi \log p\left(X; \{\Pi_k, \mu_k, \Sigma_k\}\right) = \lambda \, \nabla_\Pi \sum_{k=1}^{k} \Pi_k$$

$$\implies \sum_{i=1}^{n} \frac{\mathcal{N}(x_i; \mu_k, \Sigma_k)}{\sum_{k'=1}^{K} \Pi_{k'} \, \mathcal{N}(x_i; \mu_{k'}, \Sigma_{k'})} = \lambda \qquad k = 1 \cdots k$$

Multiplying both sides by $\Pi_k$ and summing over $k$ gives

This same algorithm emerges as MLE.

$$\lambda = \Pi_k \lambda = \sum_{k=1}^{K} \Pi_k \lambda =$$

$$\gamma_{ik} = \sum_{i=1}^{n} \gamma_{ik} = \sum_{k=1}^{K} \sum_{i=1}^{n} \gamma_{ik} = n = \sum_{i=1}^{n} \gamma_{ik} = n_k$$

$\Uparrow$

# The EM Algorithm

Algorithm for estimating GMM params is example of EM (Expectation-Maximization) Algorithm.

## EM

$X$ = observable (incomplete) data

Data $(X, Y)$    $Y$ = unobservable data

$\underbrace{\phantom{(X,Y)}}$ complete data.

## Ex

$$X \sim q(x) = \sum_{k=1}^{K} \pi_k N(x; \mu_k, \Sigma_k) = GMM$$

$X$ = observed data

$Y$ = which Gaussian $\in \{1 \cdots K\}$ generated $X$.

EM assumes that ~~q(x|θ) we want to estimate~~

① Have $X \sim q(x|\theta) = \sum_Y q(x,y|\theta)$

② Want to estimate $\theta$ by MLE.

③ Given complete data $(X, Y)$ easy to ~~xxx~~ compute MLE.

Given $\underline{X}$ and current $\theta$ est, $\theta^{old}$

$$q(Y|X, \theta^{old}) \text{ is dist on unobserved given observed}$$

$$Q(\theta, \theta^{old}) = E_{\theta^{old}} \log q(X, Y|\theta) = \sum_Y \log q(X, Y=y | \theta) \, q(Y=y | X, \theta^{old})$$

$$\theta^{new} = \arg\max_\theta Q(\theta, \theta^{old})$$

EM Produces sequence $\theta^1, \theta^2, \ldots$

Can show

① $q(X | \theta^i) \leq q(X | \theta^{i+1}) \leftarrow$ --- with equality iff $\theta^i$ is local max.

② Update eqns often tractable.

EM Ex (GMM Revisited)

① $y \in \{1, \dots K\}$    $p(y=k) = \pi_k$    $p(x) = \sum_{k=1}^{u} \pi_k \, N(x; \mu_k, \Sigma_k)$

② $x \sim N(\mu_k, \Sigma_k)$

$$X = \begin{pmatrix} -- x_1 -- \\ -- x_2 -- \\ \vdots \\ -- x_n -- \end{pmatrix} \quad Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

observed          unobserved

Have random sample
$(x_1, y_1) \dots (x_n, y_n)$
but don't observe $y$'s.

$$G = \left\{ \pi_k, \mu_k, \Sigma_k \right\}_{k=1}^{K}$$

$$\log p(X, Y | G) = \log \prod_{i=1}^{n} p(x_i, y_i | G) = \sum_{i=1}^{n} \log p(x_i, y_i | G)$$

Given $G^{old}$

$$p(y_i = k | x_i, G^{old}) = \frac{\pi_k \, N(x_i; \mu_k, \Sigma_k)}{\sum_{k'} \pi_{k'} \, N(x_i, \mu_{k'}, \Sigma_{k'})} = \gamma_{ik}$$

$$E_{G^{old}} \log p(X, Y | G) = E_{G^{old}} \sum_{i=1}^{n} \log p(x_i, y_i | G)$$

$$= \sum_{i=1}^{n} \sum_{k=1}^{u} \log p(x_i, y_i = k | G) \, p(y_i = k | x_i, G^{old})$$

$$= \sum_{i=1}^{n} \sum_{k=1}^{u} \log p(x_i, y_i = k | G) \, \gamma_{ik}$$

$$\left\{ \pi_k^{new}, \mu_k^{new}, \Sigma_k^{new} \right\}_{k=1}^{K} = \underset{\{\pi_k, \mu_k, \Sigma_k\}}{\arg \max} \sum_{i=1}^{n} \sum_{k=1}^{u} \log \pi_k \, N(x_i; \mu_k, \Sigma_k) \, \gamma_{ik}$$

$$= \underset{\{\pi_k, \mu_k, \Sigma_k\}}{\arg \max} \sum_{i=1}^{n} \sum_{k=1}^{u} \left( \log \pi_k + \log N(x_i; \mu_k, \Sigma_k) \right) \gamma_{ik}$$

Note   Can max indep over $\{\pi_k\}_{k=1}^{K}$ and $\{\mu_k, \Sigma_k\}_{k=1}^{K}$

① $\Pi^{new} = \underset{\Pi \atop \Sigma\Pi_k=1}{\arg\max} \sum_{i=1}^{n} \sum_{k=1}^{K} \log \Pi_k \gamma_{ik} = \underset{\Pi \atop \Sigma\Pi_k=1}{\arg\max} \sum_{k=1}^{K} \log \Pi_k \, n_k$

~~By Lagrange set $\nabla_\Pi \sum_{i=1}^{n}\sum_{k=1}^{K} \log \Pi_k \gamma_{ik} = \lambda \nabla_\Pi \sum \Pi_k n$~~

By Lagrange set $\quad \nabla_\Pi \sum_{k=1}^{K} \log \Pi_k \, n_k = \lambda \, \nabla_\Pi \sum_{k=1}^{K} \Pi_k$

$\Rightarrow \qquad \dfrac{n_k}{\Pi_k} = \lambda \, \mathbf{1} \qquad \Longleftrightarrow \qquad$ ~~$\frac{n_k}{\Pi_k}$~~ $= \sum_k n_k = \sum_k \lambda \Pi_k$

$$\overline{\qquad\qquad\qquad\qquad} \\ n = \lambda$$

$\Rightarrow \Pi_k = \dfrac{n_k}{n}$

② $\mu^{new} = \underset{\mu}{\arg\max} \sum_{i=1}^{n} \sum_{k=1}^{K} \log \mathcal{N}(x_i; \mu_k, \Sigma_k) \gamma_{ik}$

$\qquad = \underset{\mu}{\arg\max} \sum_{k=1}^{K} \sum_{i=1}^{n} \log \mathcal{N}(x_i; \mu_k, \Sigma_k) \gamma_{ik}$

Can max inner sum for each $\mu_k \Rightarrow$

$\mu_k^{new} = \underset{\mu_k}{\arg\max} \sum_{i=1}^{n} \log \mathcal{N}(x_i; \mu_k, \Sigma_k) \gamma_{ik}$

Setting $\nabla_{\mu_k} = 0$ gives

$0 = \sum_{i=1}^{n} \dfrac{\mathcal{N}(x_i; \mu_k, \Sigma_k) \overset{1}{\cancel{\Sigma_k^{-1}}} \Sigma_k^{-1}(x_i - \mu_k) \gamma_{ik}}{\cancel{\mathcal{N}(x_i; \mu_k, \Sigma_k)}}$

$= \cancel{\Sigma_k^{-1}} \sum_{i=1}^{n} (x_i - \mu_k) \gamma_{ik}$

$\Rightarrow \sum_{i=1}^{n} x_i \gamma_{ik} = \sum_{i=1}^{n} \mu_k \gamma_{ik} = \mu_k n_k \Rightarrow \mu_k^{new} = \dfrac{1}{n_k} \sum_{i=1}^{n} x_i \gamma_{ik}$

em_faithful. ✓

③ Can show ~~$\mu_k^{new}$~~ $\Sigma_k^{new} = \dfrac{1}{n_k} \sum_{i=1}^{n} (x_i - \mu_k^{new})(x_i - \mu_k^{new})^T \gamma_{ik}$

Do em_faithful. ✓

# Multiclass Logistic Regression

Softmax

$$q(k \mid x) = \frac{e^{w_k^t \phi(x)}}{\sum_{k'} e^{w_{k'}^t \cdot \phi(x)}} = \text{softmax fn.}$$

Many classification problems lead to posteriors having softmax formulation

Ex    Consider generative model

$$P(C = k) = \pi_k \qquad k = 1 \cdots k$$

$$X \mid C = k \sim \text{Binomial}(n, q_k)$$

$$q(C = k \mid x) = \frac{\pi_k \binom{n}{x} q_k^x (1-q_k)^{n-x}}{\sum_{k'} \pi_{k'} \binom{n}{x} q_{k'}^x (1-q_{k'})^{n-x}} = \frac{f(k,x)}{\sum_{k'} f(k',x)}$$

$$\frac{\left(\frac{q_k}{1-q_k}\right)^x \left(\pi_k (1-q_k)\right)^n}{\sum_{k'}}$$

where $w_k = \begin{pmatrix} \log \frac{q_k}{1-q_k} \\ \log \pi_k (1-q_k)^n \end{pmatrix}$

$$f(k,x) = \pi_k \, q_k^x (1-q_k)^{n-x}$$

$$= \pi_k \left(\frac{q_k}{1-q_k}\right)^x (1-q_k)^n$$

$$\phi(x) = \begin{pmatrix} x \\ 1 \end{pmatrix}$$

$$= e^{\log \pi_k \left(\frac{q_k}{1-q_k}\right)^x (1-q_k)^n}$$

$$= e^{x \log \frac{q_k}{1-q_k} + 1 \cdot \log \pi_k (1-q_k)^n} = e^{w_k^t \phi(x)}$$

$$e^{\displaystyle \frac{w_k^t \phi(x)}{\sum_{k'} e^{w_{k'}^t \phi(x)}}} = q(C = k \mid x) =$$

# Multiclass Logistic

Have feature vector $\emptyset(x)$ to be classified in $k$ class $\{1, \ldots k\}$.

## Model

$$q_k \doteq P(\text{class} = k \mid \emptyset(x)) = \frac{e^{v_k^t \emptyset(x)}}{\sum_{k'} e^{v_{k'}^t \emptyset(x)}} \quad \underset{a_k = w_k^t \emptyset(x)}{=} \quad \frac{e^{a_k}}{\sum_{k'} e^{a_{k'}}}$$

$$\frac{\partial q_k}{\partial a_j} = \frac{\left(\sum_{k'} e^{a_{k'}}\right) e^{a_k} I_{kj} - e^{a_k} e^{a_j}}{\left(\sum_{k'} e^{a_{k'}}\right)^2} = q_k I_{kj} - q_k q_j$$

$$= q_k (I_{kj} - q_j)$$

Have training data $(\emptyset(x_1), c_1), (\emptyset(x_2), c_2) \cdots (\emptyset(x_n), c_n)$

$c_i \in \{1 \cdots k\}$ are classes

$$X = \begin{pmatrix} -- & \emptyset(x_1) & -- \\ -- & \emptyset(x_2) & -- \\ \vdots & \vdots & \vdots \\ -- & \emptyset(x_n) & -- \end{pmatrix} \qquad T = \begin{pmatrix} t_{11} & \cdots & t_{1k} \\ t_{21} & \cdots & t_{2k} \\ \vdots & & \vdots \\ t_{n1} & -- & t_{nk} \end{pmatrix}$$

$$t_{ik} = \begin{cases} 1 & c_i = k \\ 0 & o.w. \end{cases}$$

$$c_i = 1 \implies t_i = 100 \cdots 0$$
$$c_i = 2 \implies t_i = 010 \cdots 0$$

$$P(T \mid w_1 \cdots w_n) = \prod_{i=1}^{n} \prod_{k=1}^{k} q_{ik}^{t_{ik}} \qquad q_{ik} = \frac{e^{w_k^t \emptyset(x_i)}}{\sum_{k'} e^{w_{k'}^t \emptyset(x_i)}}$$

$$\text{Error fn} = E(w_1 \cdots w_n) = -\log P(T \mid w_1 \cdots w_n) = -\sum_{i=1}^{n} \sum_{k=1}^{k} t_{ik} \log q_{ik}$$

$$\nabla_{w_j} E(w_1 \cdots w_n) = -\sum_{i=1}^{n} \sum_{k=1}^{k} \frac{t_{ik}}{q_{ik}} \frac{\partial q_{ik}}{\partial a_j} \nabla_{w_j} a_{ij}$$

$$= -\sum_{i=1}^{n} \sum_{k=1}^{k} \frac{t_{ik}}{q_{ik}} q_{ik} (I_{kj} - \cancel{q_{ij}} q_{ij}) \emptyset(x_i)$$

$$= -\sum_{i=1}^{n} (t_{ij} - q_{ij}) \emptyset(x_i)$$

$$= \sum_{i=1}^{n} (q_{ij} - t_{ij}) \emptyset(x_i)$$

Can show
$$\nabla_{w_k} \nabla_{w_j} E(w_1 \cdots w_n)$$
$$= \sum_{i=1}^{n} q_{ij} (I_{kj} - q_{ik}) \emptyset(x_i) \emptyset(x_i)^t$$

# Classification Trees

Two common treatments:
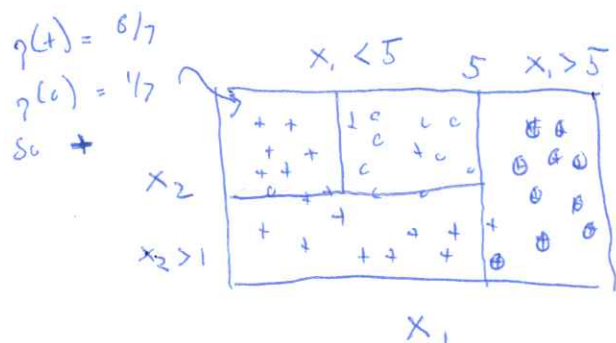
CART = Classification + Regression Trees

C4.5

We do CART

## Pictorial View

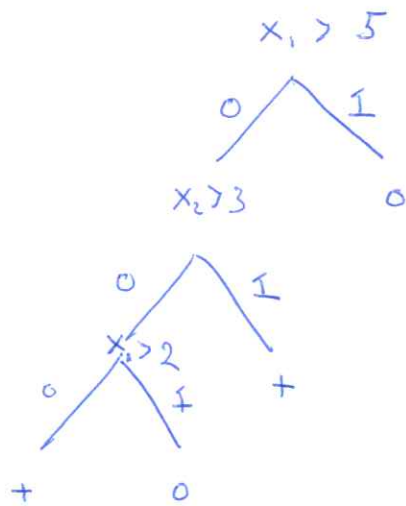Have data $(x_1, c_1), (x_2, c_2) \ldots (x_n, c_n)$

$x_i \in \mathbb{R}^2$ (2 continuous features

$c_i \in \{+, o\}$

$q(+) = 6/7$
$q(o) = 1/7$
So $+$



① Choose split $x_k < s$ dividing data into regions that are as pure as possible

② Recursively (and greedily) subdivide to increase purity (will return to stopping later)

③ Label each final region with most frequent class (or prob dist as in logistic regression)

Can visualize splits as tree

$x_1 > 5$

0 / \ 1

$x_2 > 3$                     0

0 / \ 1

$x_3 > 2$

0 / \ 1                        +

+        0

Need not be balanced!

CART works with categorical variables too

$x_k \in \{ blue, brown, hazel, green \}$

$x_k = blue$

0 / \ 1

$x_k \in \{ blue, brown \}$

0 / \ 1

Always have binary splits in CART

Do tree_prostate_cancer.r

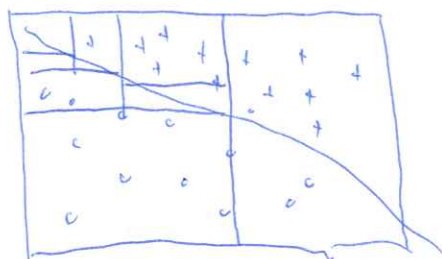# Virtues of Classification Trees

① No parametric assumptions! (So general)

② Computationally Efficient (Many variables, Large data)

③ Performs variable selection on-line

## Some Drawbacks

① Some datasets awkward

a)



Splits of form $x_k > s$

b) Binary features $x_1 \cdots x_k$

$$c(x_1 \cdots x_n) = \left( \sum_{i=1}^{k} x_i \right) \bmod 2$$

Can't classify without all features so CART fails here

② Prone to overfitting, though can be addressed.

# Choosing Splits

Splits chosen to maximize purity $\iff$ minimize impurity

## Common Impurity Measures

### ① Entropy

Have dist. $q_1 \cdots q_k$ for prob (or proportion) of classes $1 \cdots k$

Eg $\{+ \ o \ o \ + \ o \ o\}$ $q(+) = \frac{1}{3}$ ; $q(o) = \frac{2}{3}$

$$H(q) = -\sum_{k=1}^{k} q_k \log_2 q_k$$

Can show

a) $H(q) \geqslant 0$

b) $H(q) = 0 \iff q$ concentrates on $1$ class

   (Eg $q = (0,0,1,0,0)$)

c) $H(q)$ is maximal when $q_1 = q_2 = \cdots = q_k = \frac{1}{k}$

### ② Gini Index

$$G(q) = 1 - \sum_{k=1}^{k} q_k^2$$

a) $G(q) \geqslant 0$

b) $G(q) = 0 \iff q$ concentrates on $1$ class

c) $G(q)$ max when $q_1 = \cdots q_k = 1/k$

Suppose a terminal node in CART

had $(a, b, a, a, a, b, c) \implies q = \left(\frac{4}{7}, \frac{2}{7}, \frac{1}{7}\right)$

Rather than labeling as most likely class (a), we choose a $4/7$ of time, b $\frac{2}{7}$ of time

What is prob. of error.

$P(\text{correct}) = \frac{\#r}{\#r} = \frac{2}{7} \cdot \frac{2}{7} + \frac{1}{7} \cdot \frac{1}{7} = \sum_k q_k^2$

$\implies P(\text{error}) = 1 - \sum_k q_k^2 = G(q)$.