# Predicting Real-valued outputs: an introduction to Regression

**Andrew W. Moore**

**Professor**

**School of Computer Science**

**Carnegie Mellon University**
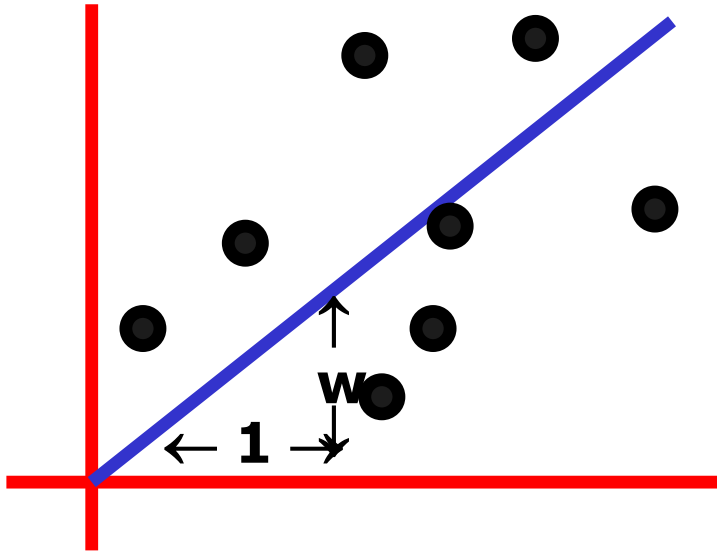
www.cs.cmu.edu/~awm

awm@cs.cmu.edu

412-268-7599

This is reordered material from the Neural Nets lecture and the "Favorite Regression Algorithms" lecture

# Single-Parameter Linear Regression

# Linear Regression

**DATASET**



| inputs | outputs |
|--------|---------|
| $x_1 = 1$ | $y_1 = 1$ |
| $x_2 = 3$ | $y_2 = 2.2$ |
| $x_3 = 2$ | $y_3 = 2$ |
| $x_4 = 1.5$ | $y_4 = 1.9$ |
| $x_5 = 4$ | $y_5 = 3.1$ |

Linear regression assumes that the expected value of the output given an input, $E[y/x]$, is linear.

Simplest case: Out($x$) = $wx$ for some unknown $w$.

Given the data, we can estimate $w$.

# 1-parameter linear regression

Assume that the data is formed by

$$y_i = wx_i + \text{noise}_i$$

where…

- the noise signals are independent
- the noise has a normal distribution with mean 0 and unknown variance $\sigma^2$

$p(y|w,x)$ has a normal distribution with

- mean $wx$
- variance $\sigma^2$

# Bayesian Linear Regression

$$p(y|w,x) = \text{Normal (mean } wx, \text{ var } \sigma^2)$$

We have a set of datapoints $(x_1, y_1)$ $(x_2, y_2)$ ... $(x_n, y_n)$ which are EVIDENCE about $w$.

We want to infer $w$ from the data.

$$p(w|x_1, x_2, x_3, \ldots x_n, y_1, y_2 \ldots y_n)$$

- You can use BAYES rule to work out a posterior distribution for $w$ given the data.

- Or you could do Maximum Likelihood Estimation

# Maximum likelihood estimation of *w*

Asks the question:

"For which value of *w* is this data most likely to have happened?"

$$<=>$$

For what *w* is

p($y_1$, $y_2$...$y_n$ | $x_1$, $x_2$, $x_3$,...$x_n$, *w*) maximized?

$$<=>$$

For what *w* is

$$\prod_{i=1}^{n} p(y_i | w, x_i) \text{ maximized}$$

For what $w$ is

$$\prod_{i=1}^{n} p(y_i | w, x_i) \text{ maximized}$$

For what $w$ is

$$\prod_{i=1}^{n} \exp\left(-\frac{1}{2}\left(\frac{y_i - wx_i}{\sigma}\right)^2\right) \text{maximized?}$$

For what $w$ is

$$\sum_{i=1}^{n} -\frac{1}{2}\left(\frac{y_i - wx_i}{\sigma}\right)^2 \text{ maximized?}$$

For what $w$ is

$$\sum_{i=1}^{n} (y_i - wx_i)^2 \text{ minimized?}$$

# Linear Regression

The maximum likelihood *w* is the one that minimizes sum-of-squares of <u>residuals</u>

$$E = \sum_i \left( y_i - w x_i \right)^2$$

$$= \sum_i y_i^2 - \left( 2 \sum x_i y_i \right) w + \left( \sum x_i^2 \right) w^2$$

We want to minimize a quadratic function of *w*.

# Linear Regression

Easy to show the sum of squares is minimized when

$$w = \frac{\sum x_i y_i}{\sum x_i^2}$$

The maximum likelihood model is

$$Out(x) = wx$$

We can use it for prediction

# Linear Regression

Easy to show the sum of squares is minimized when

$$W = \frac{\sum x_i y_i}{\sum x_i^2}$$

The maximum likelihood model is

$$Out(x) = wx$$

We can use it for prediction



p(w)

w

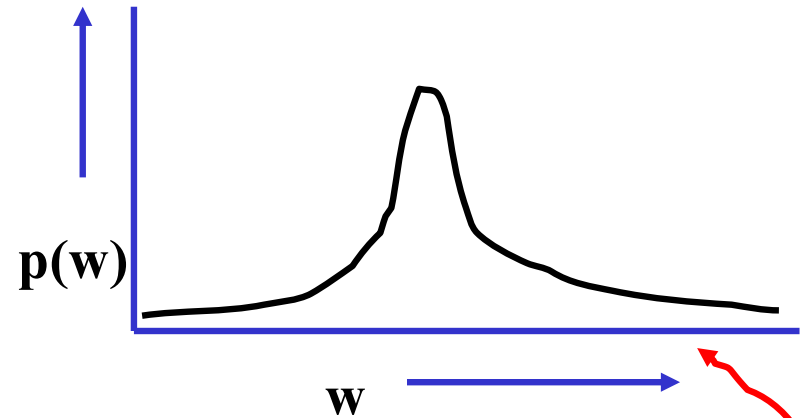**Note:** In Bayesian stats you'd have ended up with a prob dist of $w$

And predictions would have given a prob dist of expected output

Often useful to know your confidence. Max likelihood can give some kinds of confidence too.

# Multivariate Linear Regression

# Multivariate Regression

What if the inputs are vectors?

$$\begin{array}{c} 3 \\ \cdot 4 \\ 6 \cdot \\ \cdot 5 \\ \cdot 8 \\ \cdot 10 \end{array}$$

**2-d input example**

$x_2$

$x_1 \longrightarrow$

Dataset has form

| | |
|---|---|
| $\mathbf{x_1}$ | $y_1$ |
| $\mathbf{x_2}$ | $y_2$ |
| $\mathbf{x_3}$ | $y_3$ |
| $\vdots$ | $\vdots$ |
| $\mathbf{x_R}$ | $y_R$ |

# Multivariate Regression

Write matrix X and Y thus:

$$\mathbf{X} = \begin{bmatrix} ....\mathbf{X}_1.... \\ ....\mathbf{X}_2.... \\ M \\ ....\mathbf{X}_R.... \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & ... & x_{1m} \\ x_{21} & x_{22} & ... & x_{2m} \\ & & M & \\ x_{R1} & x_{R2} & ... & x_{Rm} \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ M \\ y_R \end{bmatrix}$$

(there are $R$ datapoints. Each input has $m$ components)

The linear regression model assumes a vector $\mathbf{w}$ such that

$$\text{Out}(\mathbf{x}) = \mathbf{w}^T\mathbf{x} = w_1 x[1] + w_2 x[2] + ....w_m x[D]$$

The max. likelihood $\mathbf{w}$ is $\mathbf{w} = (X^T X)^{-1}(X^T Y)$

# Multivariate Regression

Write matrix X and Y thus:

$$\mathbf{X} = \begin{bmatrix} ....\mathbf{X}_1.... \\ ....\mathbf{X}_2.... \\ M \\ ....\mathbf{X}_R.... \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & ... & x_{1m} \\ x_{21} & x_{22} & ... & x_{2m} \\ & & M & \\ x_{R1} & x_{R2} & ... & x_{Rm} \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ M \\ y_R \end{bmatrix}$$

(there are $R$ datapoints. Each input

**IMPORTANT EXERCISE: PROVE IT !!!!!**

The linear regression model assumes a vector $\mathbf{w}$ such that

$$\text{Out}(\mathbf{x}) = \mathbf{w}^T\mathbf{x} = w_1x[1] + w_2x[2] + ....w_mx[D]$$

The max. likelihood $\mathbf{w}$ is $\mathbf{w} = (X^TX)^{-1}(X^TY)$

# Multivariate Regression (con't)

The max. likelihood $w$ is $w = (X^TX)^{-1}(X^TY)$

$X^TX$ is an $m \times m$ matrix:  i,j'th elt is $\displaystyle\sum_{k=1}^{R} x_{ki}x_{kj}$

$X^TY$ is an $m$-element vector:  i'th elt $\displaystyle\sum_{k=1}^{R} x_{ki}y_{k}$

# Constant Term in Linear Regression

# What about a constant term?

We may expect linear data that does not go through the origin.

Statisticians and Neural Net Folks all agree on a simple obvious hack.

Can you guess??

# The constant term

- The trick is to create a fake input "$X_0$" that always takes the value 1

| $X_1$ | $X_2$ | $Y$ |
|-------|-------|-----|
| 2 | 4 | 16 |
| 3 | 4 | 17 |
| 5 | 5 | 20 |

Before:

$Y = w_1 X_1 + w_2 X_2$

...has to be a poor model

| $X_0$ | $X_1$ | $X_2$ | $Y$ |
|-------|-------|-------|-----|
| 1 | 2 | 4 | 16 |
| 1 | 3 | 4 | 17 |
| 1 | 5 | 5 | 20 |

After:

$Y = w_0 X_0 + w_1 X_1 + w_2 X_2$

$= w_0 + w_1 X_1 + w_2 X_2$

...has a fine constant term

In this example, You should be able to see the MLE $w_0$, $w_1$ and $w_2$ by inspection

# Non-linear Regression

# Non-linear Regression

- Suppose you know that y is related to a function of x in such a way that the predicted values have a non-linear dependence on w, e.g:

| $x_i$ | $y_i$ |
|-------|-------|
| ½ | ½ |
| 1 | 2.5 |
| 2 | 3 |
| 3 | 2 |
| 3 | 3 |



Assume $y_i \sim N(\sqrt{w + x_i}, \sigma^2)$

What's the MLE estimate of w?

# Non-linear MLE estimation

$$\underset{w}{\text{argmax}} \log p(y_1, y_2, \ldots, y_R \mid x_1, x_2, \ldots, x_R, \sigma, w) =$$

$$\underset{w}{\text{argmin}} \sum_{i=1}^{R} \left( y_i - \sqrt{w + x_i} \right)^2 =$$

Assuming i.i.d. and then plugging in equation for Gaussian and simplifying.

$$\left( w \text{ such that } \sum_{i=1}^{R} \frac{y_i - \sqrt{w + x_i}}{\sqrt{w + x_i}} = 0 \right) =$$

Setting dLL/dw equal to zero

# Non-linear MLE estimation

$$\underset{w}{\operatorname{argmax}} \log p(y_1, y_2, \ldots, y_R \mid x_1, x_2, \ldots, x_R, \sigma, w) =$$

$$\underset{w}{\operatorname{argmin}} \sum_{i=1}^{R} \left( y_i - \sqrt{w + x_i} \right)^2 =$$

Assuming i.i.d. and then plugging in equation for Gaussian and simplifying.

$$\left( w \text{ such that } \sum_{i=1}^{R} \frac{y_i - \sqrt{w + x_i}}{\sqrt{w + x_i}} = 0 \right) =$$

Setting dLL/dw equal to zero

We're down the algebraic toilet

So guess what we do?

# Non-linear MLE estimation

$$\text{argmax}_w \log p(y_1, y_2, \ldots, y_R \mid x_1, x_2, \ldots, x_R, \sigma, w) =$$

$$\left( \ldots \cdot w + x_i \right)^2 =$$

Assuming i.i.d. and then plugging in equation for Gaussian and simplifying.

$$\frac{\ldots + x_i}{\ldots} = 0 \right) =$$

Setting dLL/dw equal to zero

Common (but not only) approach:
Numerical Solutions:
- Line Search
- Simulated Annealing
- Gradient Descent
- Conjugate Gradient
- Levenberg Marquart
- Newton's Method

*Also, special purpose statistical-optimization-specific tricks such as E.M. (See Gaussian Mixtures lecture for introduction)*

We're down the algebraic toilet

So guess what we do?

23

# Polynomial Regression

24

# Polynomial Regression

So far we've mainly been dealing with linear regression

| $X_1$ | $X_2$ | $Y$ |
|-------|-------|-----|
| 3 | 2 | 7 |
| 1 | 1 | 3 |
| : | : | : |

$$\mathbf{X}= \begin{array}{|c|c|} \hline 3 & 2 \\ \hline 1 & 1 \\ \hline : & : \\ \hline \end{array} \qquad \mathbf{y}= \begin{array}{|c|} \hline 7 \\ \hline 3 \\ \hline : \\ \hline \end{array}$$

$z_1=(3,2)..$      $y_1=7..$

$$\mathbf{Z}= \begin{array}{|c|c|c|} \hline 1 & 3 & 2 \\ \hline 1 & 1 & 1 \\ \hline : & & : \\ \hline \end{array} \qquad \mathbf{y}= \begin{array}{|c|} \hline 7 \\ \hline 3 \\ \hline : \\ \hline \end{array}$$

$\mathbf{z}_1=(1,3,2)..$      $y_1=7..$

$\mathbf{z}_k=(1,x_{k1},x_{k2})$

$$\beta=(\mathbf{Z}^T\mathbf{Z})^{-1}(\mathbf{Z}^T\mathbf{y})$$

$$y^{est} = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

# Quadratic Regression

It's trivial to do linear fits of fixed nonlinear basis functions

| $X_1$ | $X_2$ | $Y$ |
|-------|-------|-----|
| 3 | 2 | 7 |
| 1 | 1 | 3 |
| : | : | : |

$$\mathbf{X}= \begin{array}{|c|c|} \hline 3 & 2 \\ \hline 1 & 1 \\ \hline : & : \\ \hline \end{array} \qquad \mathbf{y}= \begin{array}{|c|} \hline 7 \\ \hline 3 \\ \hline : \\ \hline \end{array}$$

$y_1=7..$

$$\mathbf{Z}= \begin{array}{|c|c|c|c|c|c|} \hline 1 & 3 & 2 & 9 & 6 & 4 \\ \hline 1 & 1 & 1 & 1 & 1 & 1 \\ \hline : & & & & & : \\ \hline \end{array} \qquad \mathbf{y}= \begin{array}{|c|} \hline 7 \\ \hline 3 \\ \hline : \\ \hline \end{array}$$

$\mathbf{z}=(1 \ , \ x_1, \ x_2, \ x_1^2, x_1 x_2, x_2^2)$

$$\beta=(\mathbf{Z}^T\mathbf{Z})^{-1}(\mathbf{Z}^T\mathbf{y})$$

$$y^{est} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1^2 + \beta_4 x_1 x_2 + \beta_5 x_2^2$$

# Quadratic Regression

It's tri...

| $X_1$ | $X_2$ |
|-------|-------|
| 3 | 2 |
| 1 | 1 |
| : | : |

$$\mathbf{Z}= \begin{bmatrix} 1 \\ 1 \\ : \end{bmatrix}$$

$z=(1 ...$

Each component of a z vector is called a term.

Each column of the Z matrix is called a term column

How many terms in a quadratic regression with $m$ inputs?

- 1 constant term

- m linear terms

- (m+1)-choose-2 = m(m+1)/2 quadratic terms

(m+2)-choose-2 terms in total $= O(m^2)$

Note that solving $\beta = (\mathbf{Z}^T\mathbf{Z})^{-1}(\mathbf{Z}^T\mathbf{y})$ is thus $O(m^6)$

# Q$^{th}$-degree polynomial Regression

| $X_1$ | $X_2$ | $Y$ |
|-------|-------|-----|
| 3 | 2 | 7 |
| 1 | 1 | 3 |
| : | : | : |

$$\mathbf{X}= \begin{array}{|c|c|} \hline 3 & 2 \\ \hline 1 & 1 \\ \hline : & : \\ \hline \end{array} \quad \mathbf{y}= \begin{array}{|c|} \hline 7 \\ \hline 3 \\ \hline : \\ \hline \end{array}$$

$$\mathbf{z}= \begin{array}{|c|c|c|c|c|c|} \hline 1 & 3 & 2 & 9 & 6 & ... \\ \hline 1 & 1 & 1 & 1 & 1 & ... \\ \hline : & & & & & ... \\ \hline \end{array} \quad \mathbf{y}= \begin{array}{|c|} \hline 7 \\ \hline 3 \\ \hline : \\ \hline \end{array}$$

***z*=(all products of powers of inputs in which sum of powers is q or less)**

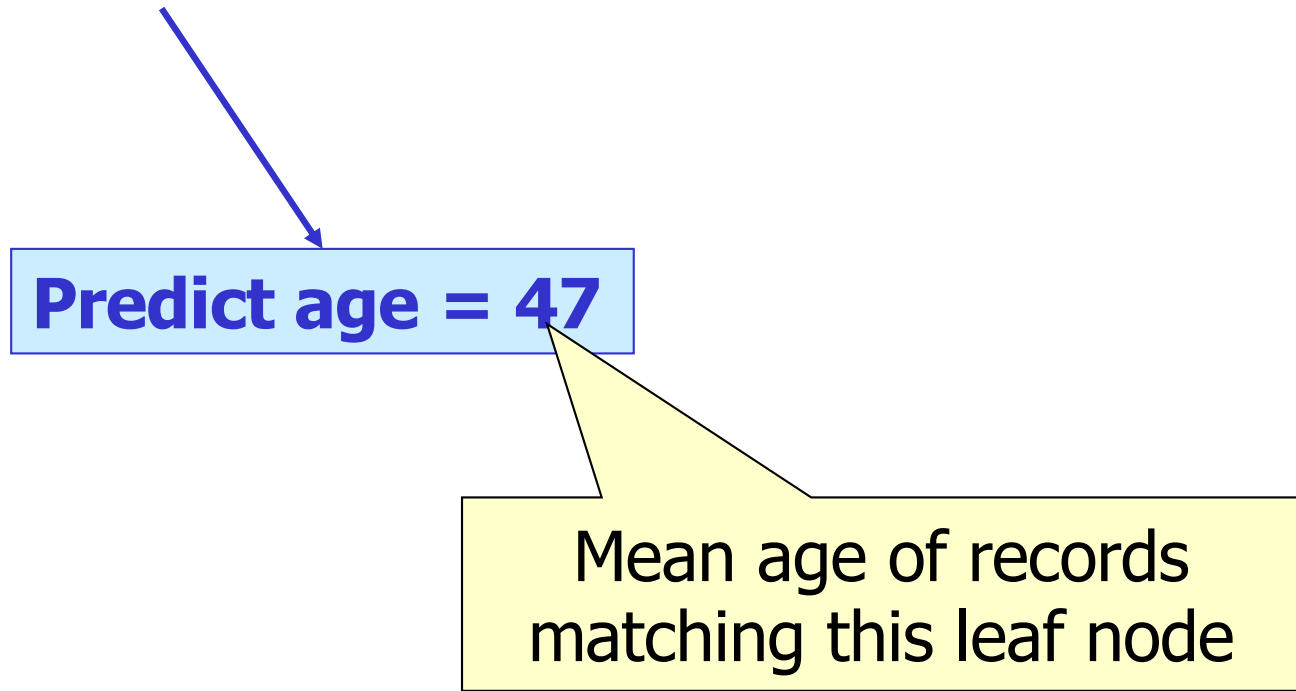$$\beta=(\mathbf{Z}^T\mathbf{Z})^{-1}(\mathbf{Z}^T\mathbf{y})$$

$$y^{est} = \beta_0 + \beta_1 x_1 + ...$$

# Regression Trees

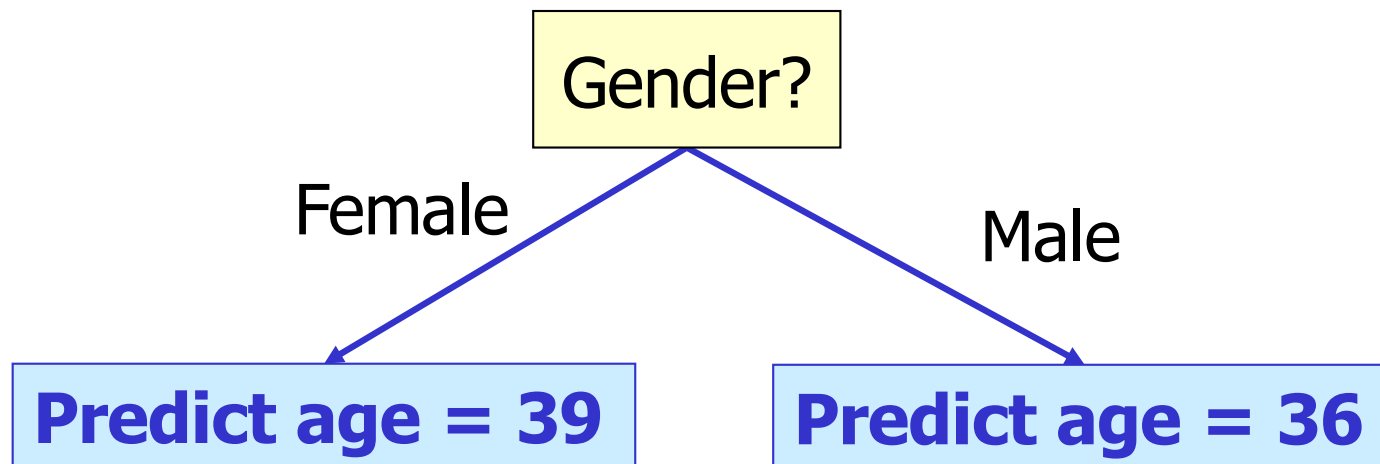# Regression Trees

- "Decision trees for regression"

# A regression tree leaf

**Predict age = 47**

Mean age of records matching this leaf node

# A one-split regression tree



Gender?

Female → **Predict age = 39**

Male → **Predict age = 36**

# Choosing the attribute to split on

| Gender | Rich? | Num. Children | Num. Beany Babies | Age |
|--------|-------|---------------|-------------------|-----|
| Female | No | 2 | 1 | 38 |
| Male | No | 0 | 0 | 24 |
| Male | Yes | 0 | 5+ | 72 |
| : | : | : | : | : |

- We can't use information gain.
- What should we use?

# Choosing the attribute to split on

| Gender | Rich? | Num. Children | Num. Beany Babies | Age |
|--------|-------|---------------|-------------------|-----|
| Female | No | 2 | 1 | 38 |
| Male | No | 0 | 0 | 24 |
| Male | Yes | 0 | 5+ | 72 |
| : | : | : | : | : |

MSE(Y|X) = The expected squared error if we must predict a record's Y value given only knowledge of the record's X value

If we're told $x=j$, the smallest expected error comes from predicting the mean of the Y-values among those records in which $x=j$. Call this mean quantity $\mu_y^{x=j}$

Then…

$$MSE(Y \mid X) = \frac{1}{R} \sum_{j=1}^{N_X} \sum_{(k \text{ such that } x_k = j)} (y_k - \mu_y^{x=j})^2$$

# Choosing the attribute to split on

| Gender | Rich? | Num. Children | Num. Beany Babies | Age |
|--------|-------|---------------|-------------------|-----|
| Female | N | | | |
| Male | N | | | |
| Male | Y | | | |
| : | : | | | |

Regression tree attribute selection: greedily choose the attribute that minimizes MSE(Y|X)

Guess what we do about real-valued inputs?

Guess how we prevent overfitting

MSE(Y|X) = The expected squared error if we must predict a record's Y value given only knowledge of the record's X value
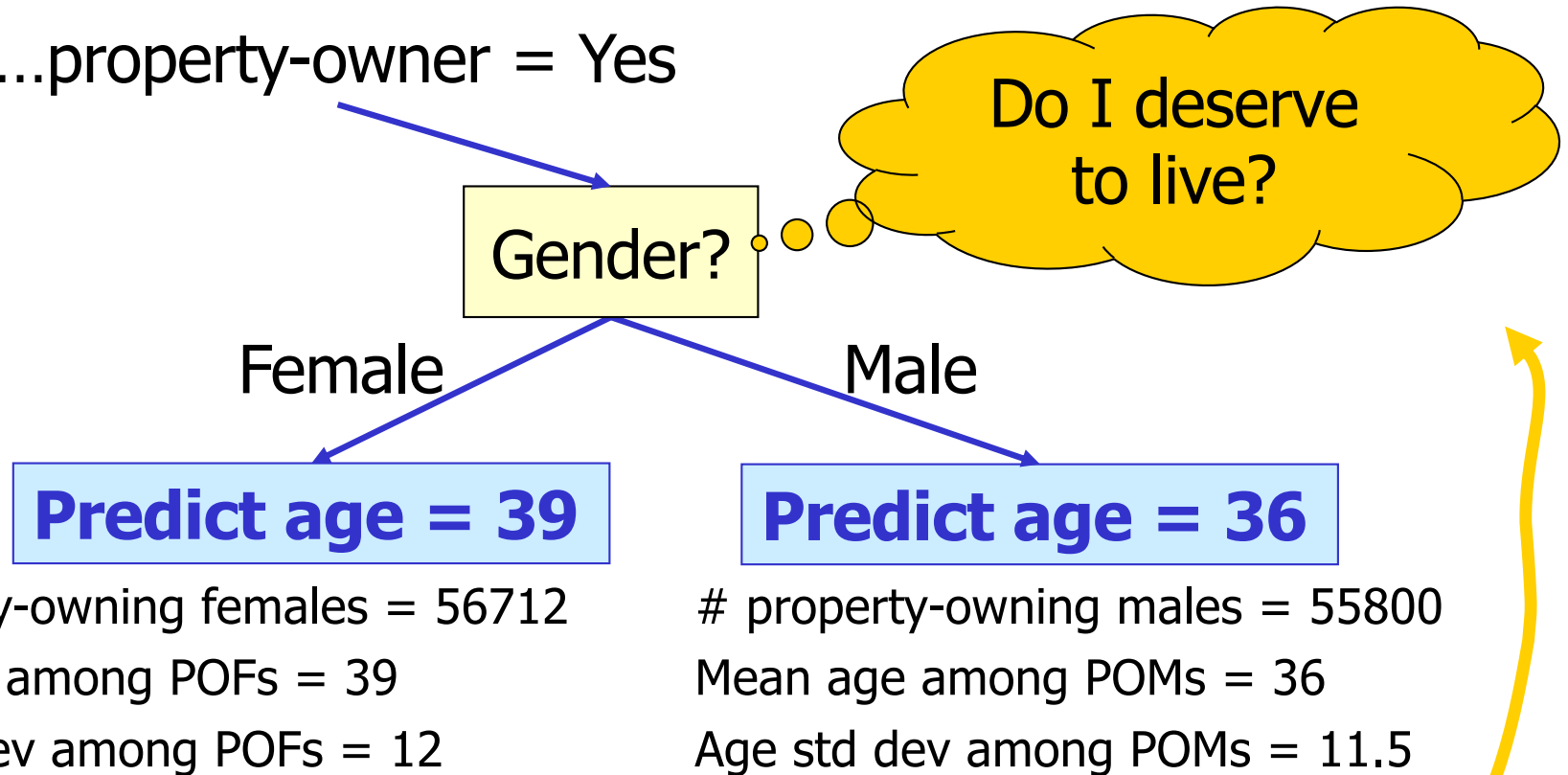
If we're told $x=j$, the smallest expected error comes from predicting the mean of the Y-values among those records in which $x=j$. Call this mean quantity $\mu_y^{x=j}$

Then...

$$MSE(Y \mid X) = \frac{1}{R} \sum_{j=1}^{N_X} \sum_{(k \, such \, that \, x_k = j)} (y_k - \mu_y^{x=j})^2$$

# Pruning Decision

...property-owner = Yes

Do I deserve to live?

Gender?

Female | Male

**Predict age = 39** | **Predict age = 36**

# property-owning females = 56712

Mean age among POFs = 39

Age std dev among POFs = 12

# property-owning males = 55800

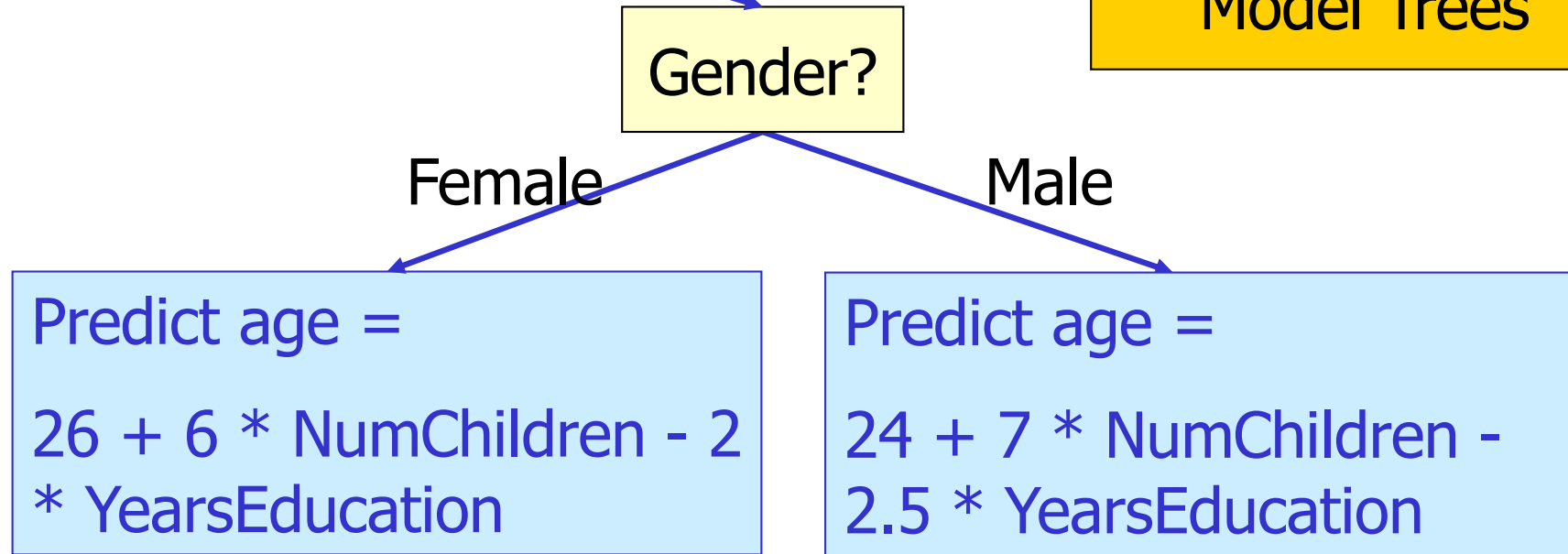Mean age among POMs = 36

Age std dev among POMs = 11.5

Use a standard Chi-squared test of the null-hypothesis "these two populations have the same mean" and Bob's your uncle.

# Linear Regression Trees

...property-owner = Yes

**Also known as "Model Trees"**

Gender?

Female        Male

Predict age =

26 + 6 * NumChildren - 2 * YearsEducation

Predict age =

24 + 7 * NumChildren - 2.5 * YearsEducation

Leaves contain linear functions (trained using linear regression on all records matching that leaf)

Split attribute chosen to minimize MSE of regressed children.

Pruning with a different Chi-squared

# Linear Regression Trees

...property-owner = Yes

Gender?

Also known as "Model Trees"
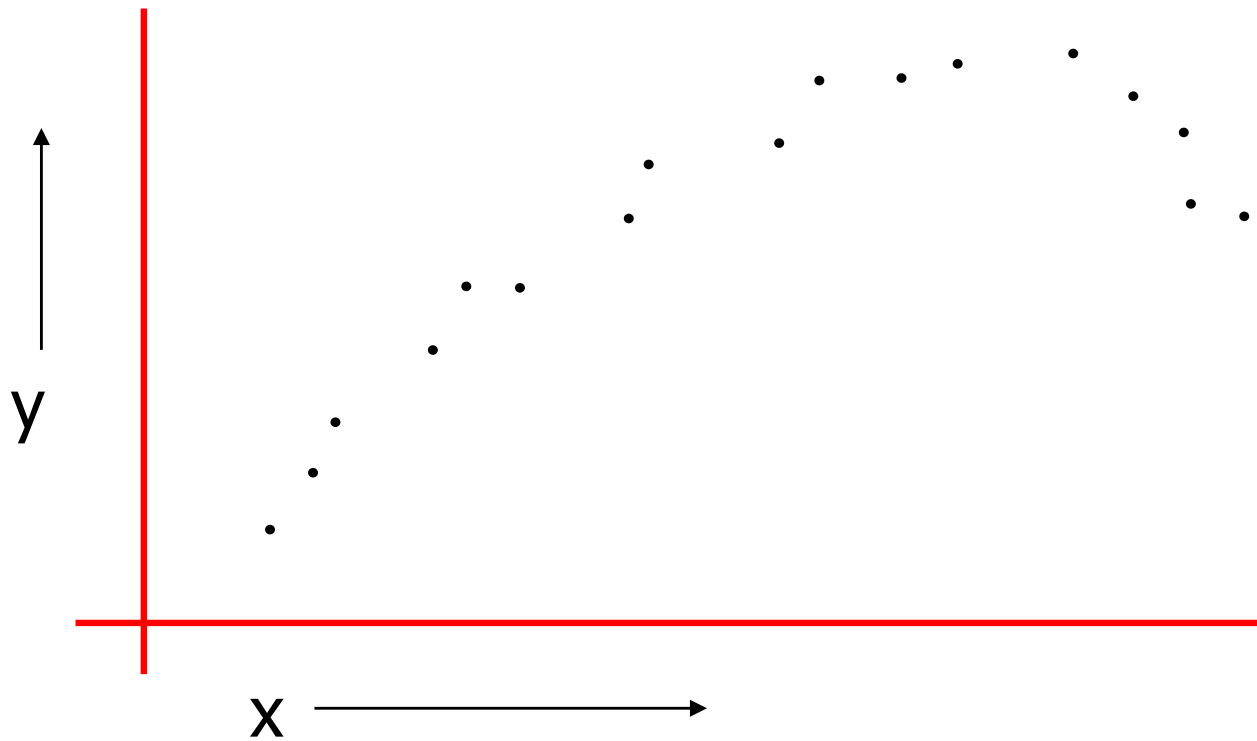
Female

Predict age =

26 + 6 * N...
* YearsE...

Detail: You typically ignore any categorical attribute that has been tested on higher up in the tree during the regression. But use all untested attributes, and use real-valued attributes even if they've been tested above

Leaves contain...
functions (traine...
linear regression...
records matching t...

...chosen to minimize
...regressed children.

Pruning with a different Chi-squared

# Test your understanding

Assuming regular regression trees, can you sketch a graph of the fitted function $y^{est}(x)$ over this diagram?

# Test your understanding

Assuming linear regression trees, can you sketch a graph of the fitted function $y^{est}(x)$ over this diagram?