

Note to other teachers and users of these slides. Andrew would be delighted if you found this source material useful in giving your own lectures. Feel free to use these slides verbatim, or to modify them to fit your own needs. PowerPoint originals are available. If you make use of a significant portion of these slides in your own lecture, please include this message, or the following link to the source repository of Andrew's tutorials: <http://www.cs.cmu.edu/~awm/tutorials> . Comments and corrections gratefully received.

Learning Gaussian Bayes Classifiers

Andrew W. Moore
Associate Professor
School of Computer Science
Carnegie Mellon University

www.cs.cmu.edu/~awm

awm@cs.cmu.edu

412-268-7599

Maximum Likelihood learning of Gaussians for Classification

- Why we should care
- 3 seconds to teach you a new learning algorithm
- What if there are 10,000 dimensions?
- What if there are categorical inputs?
- Examples “out the wazoo”

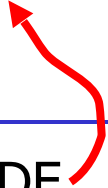
Why we should care

- One of the original “Data Mining” algorithms
- Very simple and effective
- Demonstrates the usefulness of our earlier groundwork

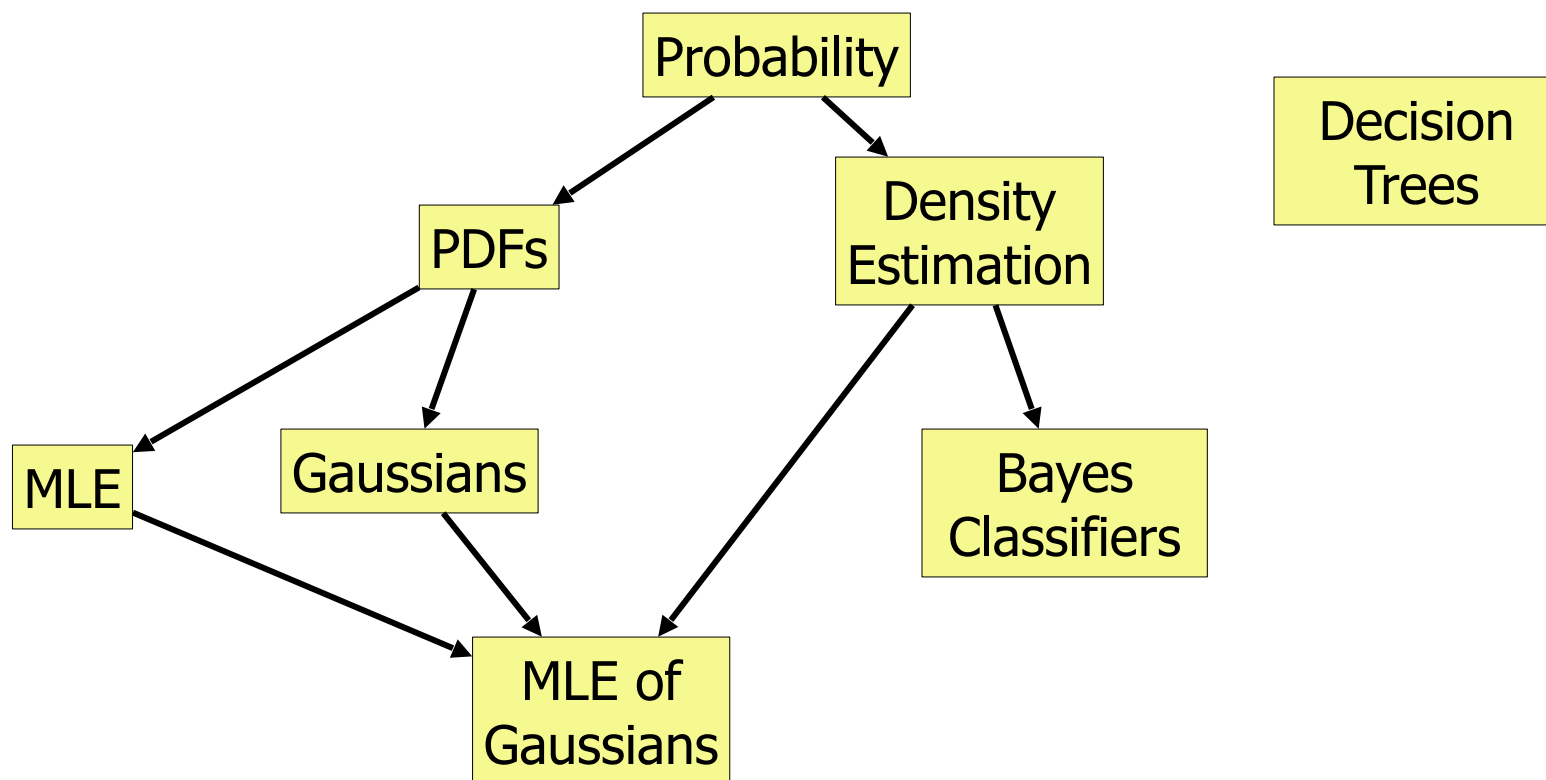
Where we were at the end of the MLE lecture...

		Categorical inputs only	Real-valued inputs only	Mixed Real / Cat okay
<div>Inputs</div> <div><div>Classifier</div><div>Predict category</div></div>	Joint BC Naïve BC		Dec Tree	
<div>Inputs</div> <div><div>Density Estimator</div><div>Prob- ability</div></div>	Joint DE Naïve DE	Gauss DE		
<div>Inputs</div> <div><div>Regressor</div><div>Predict real no.</div></div>				

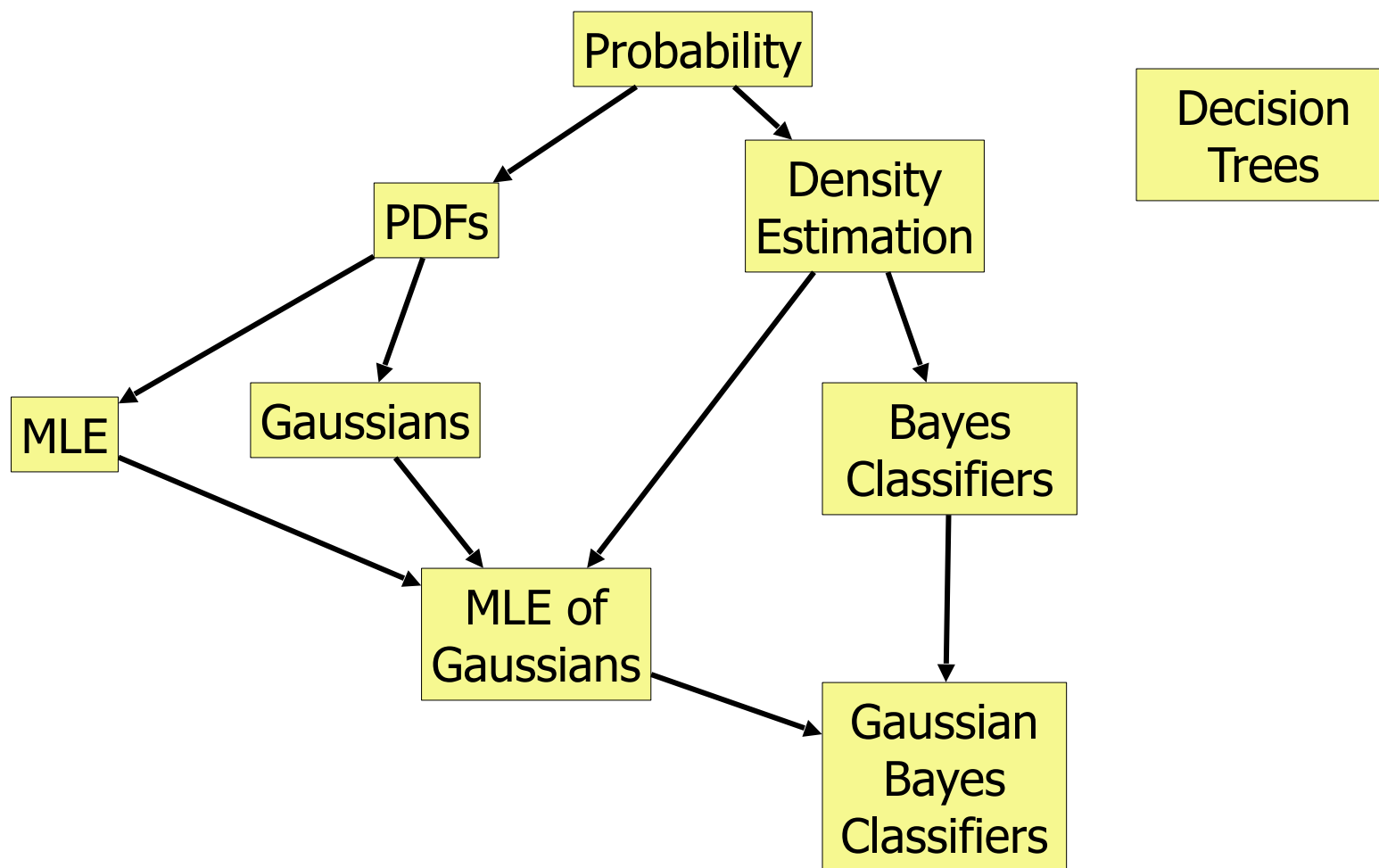
This lecture...

			Categorical inputs only	Real-valued inputs only	Mixed Real / Cat okay
Inputs — — — — —	Classifier	Predict category	Joint BC Naïve BC	Gauss BC 	Dec Tree
Inputs — — — — —	Density Estimator	Prob- ability	Joint DE Naïve DE	Gauss DE	
Inputs — — — — —	Regressor	Predict real no.			

Road Map



Road Map



Gaussian Bayes Classifier Assumption

- The i 'th record in the database is created using the following algorithm
 1. Generate the output (the "class") by drawing $y_i \sim \text{Multinomial}(p_1, p_2, \dots, p_{N_y})$
 2. Generate the inputs from a Gaussian PDF that depends on the value of y_i :

$$\mathbf{x}_i \sim N(\mu_i, \Sigma_i).$$

Test your understanding. Given N_y classes and m input attributes, how many distinct scalar parameters need to be estimated?

MLE Gaussian Bayes Classifier

Let DB_i = Subset of database DB in which the output class is $y = i$

$$p_i^{mle} = \frac{|DB_i|}{|DB|}$$

1. Generate the output (the "class") by drawing $y_i \sim \text{Multinomial}(p_1, p_2, \dots, p_{N_y})$
2. Generate the inputs from a Gaussian PDF that depends on the value of y_i :

$$\mathbf{x}_i \sim N(\mu_i, \Sigma_i).$$

Test your understanding. Given N_y classes and m input attributes, how many distinct scalar parameters need to be estimated?

MLE Gaussian Bayes Classifier

Let DB_i = Subset of database DB in which the output class is $y = i$

the database is created
g algorithm

$(\mu_i^{\text{mle}}, \Sigma_i^{\text{mle}}) = \text{MLE Gaussian for } DB_i$

2. Generate the inputs from Gaussian PDF that depends on the value of y_i :

$$\mathbf{x}_i \sim N(\mu_i, \Sigma_i).$$

Test your understanding. Given N_y classes and m input attributes, how many distinct scalar parameters need to be estimated?

MLE Gaussian Bayes Classifier

Let DB_i = Subset of database DB in which the output class is $y = i$

the database is created
g algorithm

($\mu_i^{mle}, \Sigma_i^{mle}$) = MLE Gaussian for DB_i

2. Generate the inputs from Gaussian PDF that depends on the value of y_i :

$$\mathbf{x}_i \sim N(\mu_i, \Sigma_i).$$

$$\mu_i^{mle} = \frac{1}{|DB_i|} \sum_{\mathbf{x}_k \in DB_i} \mathbf{x}_k$$

$$\Sigma_i^{mle} = \frac{1}{|DB_i|} \sum_{\mathbf{x}_k \in DB_i} (\mathbf{x}_k - \mu_i^{mle})(\mathbf{x}_k - \mu_i^{mle})^T$$

distinct scalar parameters need to be estimated?

Gaussian Bayes Classification

$$P(y = i \mid \mathbf{x}) = \frac{p(\mathbf{x} \mid y = i)P(y = i)}{p(\mathbf{x})}$$

Gaussian Bayes Classification

$$P(y = i \mid \mathbf{x}) = \frac{p(\mathbf{x} \mid y = i)P(y = i)}{p(\mathbf{x})}$$

$$P(y = i \mid \mathbf{x}) = \frac{\frac{1}{(2\pi)^{m/2} \|\Sigma_i\|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x}_k - \boldsymbol{\mu}_i)^T \Sigma_i (\mathbf{x}_k - \boldsymbol{\mu}_i)\right] p_i}{p(\mathbf{x})}$$

How do we deal with that?

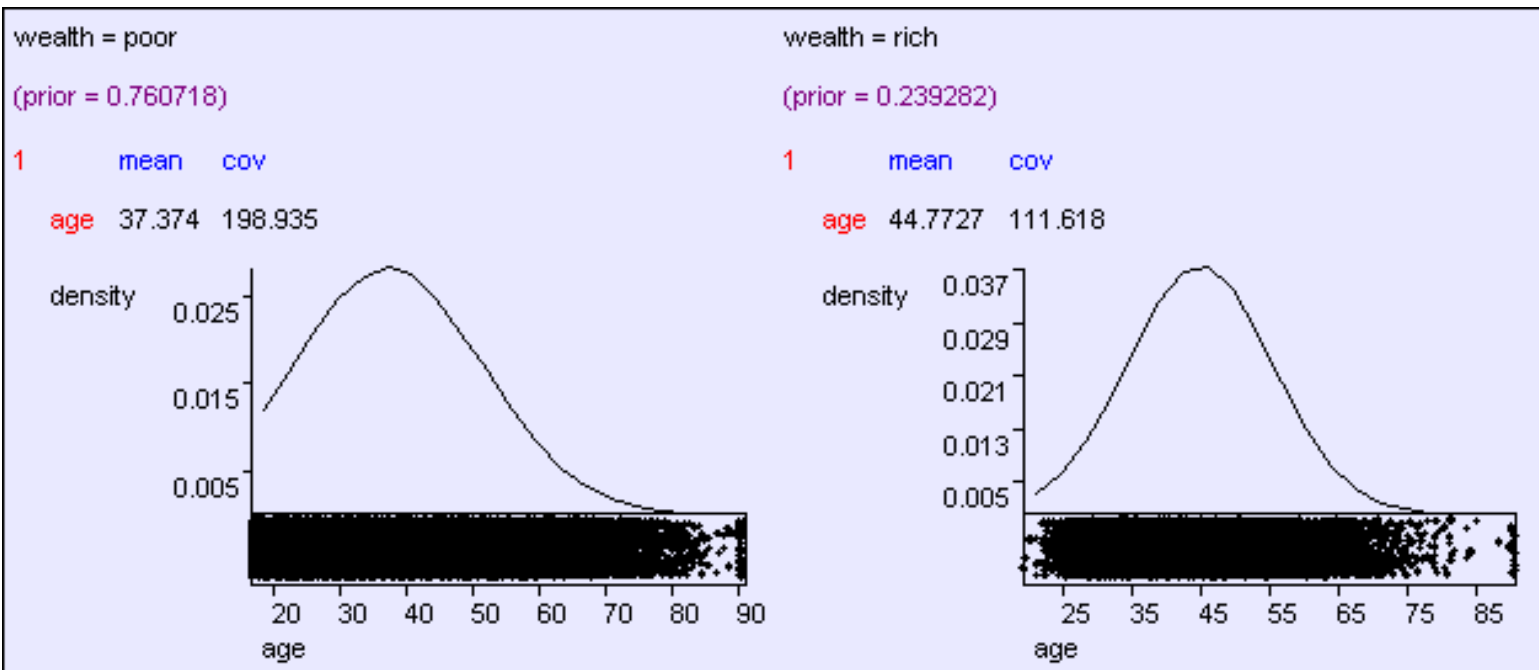


Here is a dataset

age	employe	education	edun	marital	...	job	relation	race	gender	hour	country	wealth
					...							
39	State_gov	Bachelors	13	Never_mar	...	Adm_cleric	Not_in_fan	White	Male	40	United_Stat	poor
51	Self_emp	Bachelors	13	Married	...	Exec_man	Husband	White	Male	13	United_Stat	poor
39	Private	HS_grad	9	Divorced	...	Handlers_c	Not_in_fan	White	Male	40	United_Stat	poor
54	Private	11th	7	Married	...	Handlers_c	Husband	Black	Male	40	United_Stat	poor
28	Private	Bachelors	13	Married	...	Prof_speci	Wife	Black	Female	40	Cuba	poor
38	Private	Masters	14	Married	...	Exec_man	Wife	White	Female	40	United_Stat	poor
50	Private	9th	5	Married_sp	...	Other_serv	Not_in_fan	Black	Female	16	Jamaica	poor
52	Self_emp	HS_grad	9	Married	...	Exec_man	Husband	White	Male	45	United_Stat	rich
31	Private	Masters	14	Never_mar	...	Prof_speci	Not_in_fan	White	Female	50	United_Stat	rich
42	Private	Bachelors	13	Married	...	Exec_man	Husband	White	Male	40	United_Stat	rich
37	Private	Some_coll	10	Married	...	Exec_man	Husband	Black	Male	80	United_Stat	rich
30	State_gov	Bachelors	13	Married	...	Prof_speci	Husband	Asian	Male	40	India	rich
24	Private	Bachelors	13	Never_mar	...	Adm_cleric	Own_child	White	Female	30	United_Stat	poor
33	Private	Assoc_acc	12	Never_mar	...	Sales	Not_in_fan	Black	Male	50	United_Stat	poor
41	Private	Assoc_voc	11	Married	...	Craft_repa	Husband	Asian	Male	40	*MissingV	rich
34	Private	7th_8th	4	Married	...	Transport	Husband	Amer_Indi	Male	45	Mexico	poor
26	Self_emp	HS_grad	9	Never_mar	...	Farming_fi	Own_child	White	Male	35	United_Stat	poor
33	Private	HS_grad	9	Never_mar	...	Machine_c	Unmarried	White	Male	40	United_Stat	poor
38	Private	11th	7	Married	...	Sales	Husband	White	Male	50	United_Stat	poor
44	Self_emp	Masters	14	Divorced	...	Exec_man	Unmarried	White	Female	45	United_Stat	rich
41	Private	Doctorate	16	Married	...	Prof_speci	Husband	White	Male	60	United_Stat	rich
:	:	:	:	:	:	:	:	:	:	:	:	:

48,000 records, 16 attributes [Kohavi 1995]

Predicting wealth from age



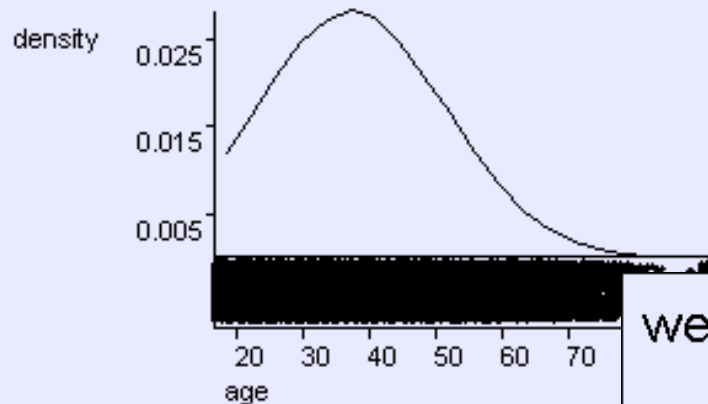
Predicting wealth from age

wealth = poor

(prior = 0.760718)

1 mean cov

age 37.374 198.935

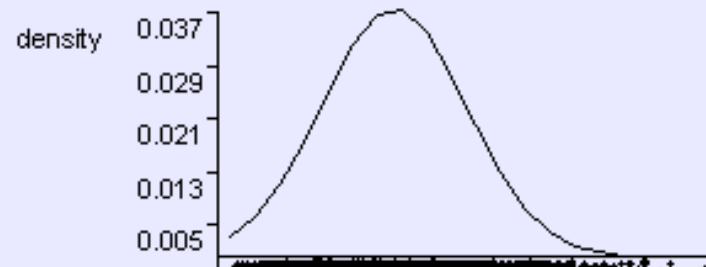


wealth = rich

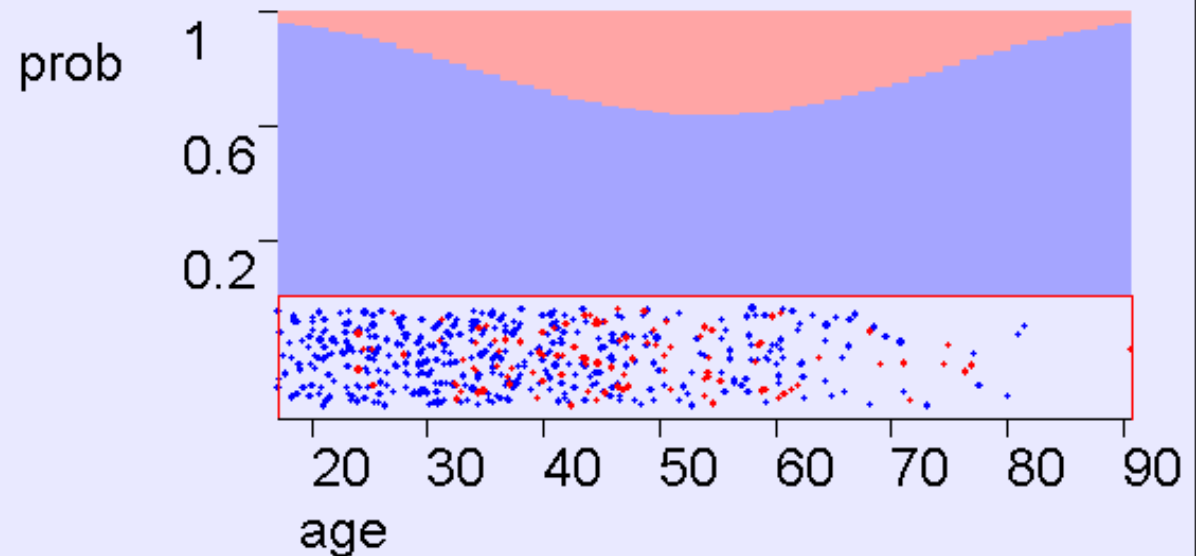
(prior = 0.239282)

1 mean cov

age 44.7727 111.618



wealth values: poor rich

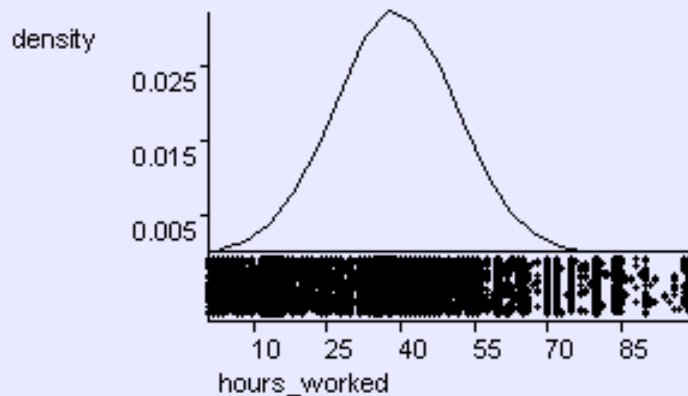


Wealth from hours worked

wealth = poor

(prior = 0.760718)

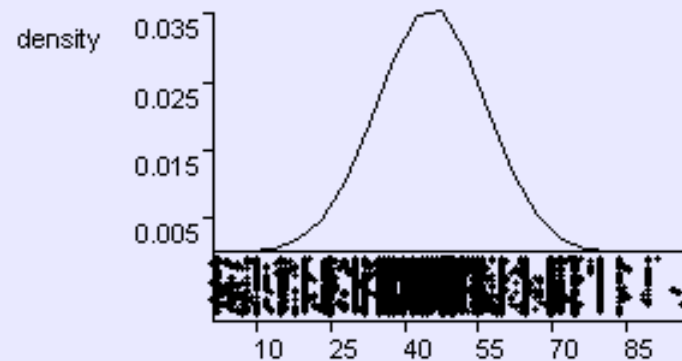
1 mean cov
hours_worked 38.84 152.692



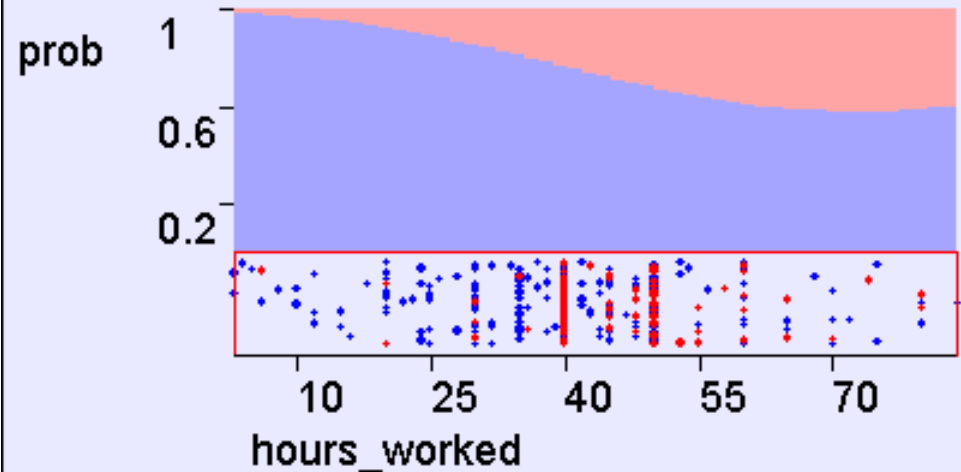
wealth = rich

(prior = 0.239282)

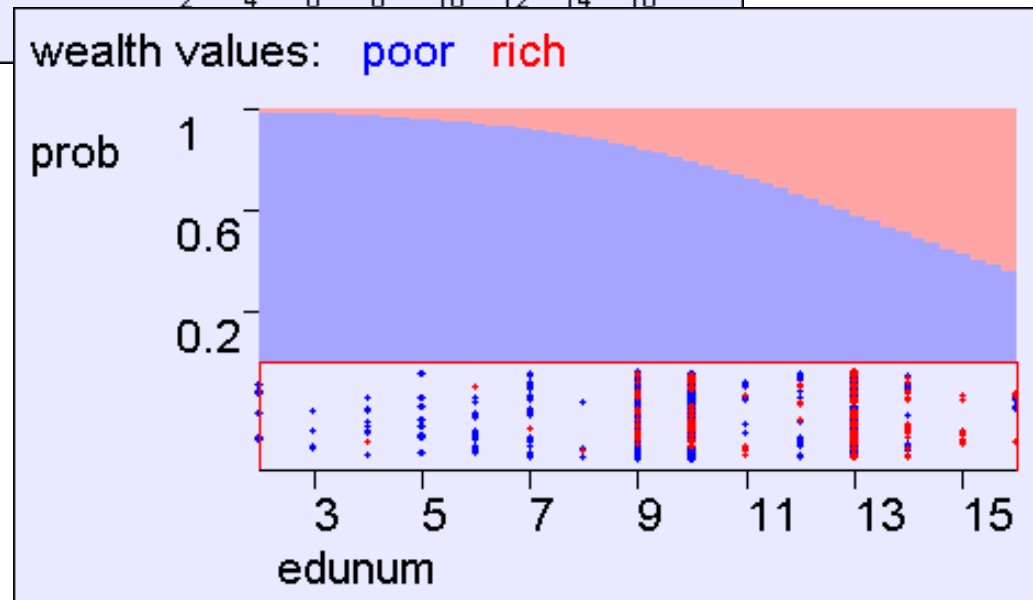
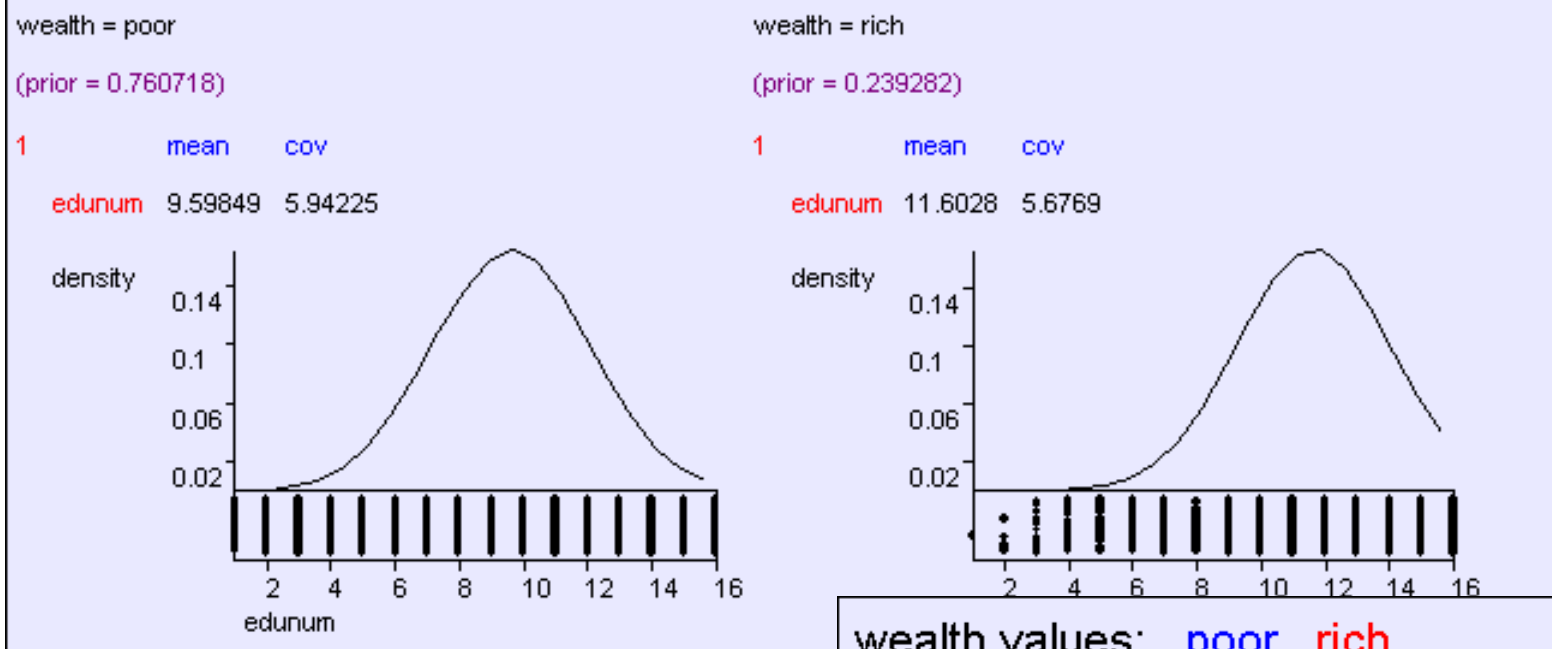
1 mean cov
hours_worked 45.4529 123.014



wealth values: poor rich



Wealth from years of education



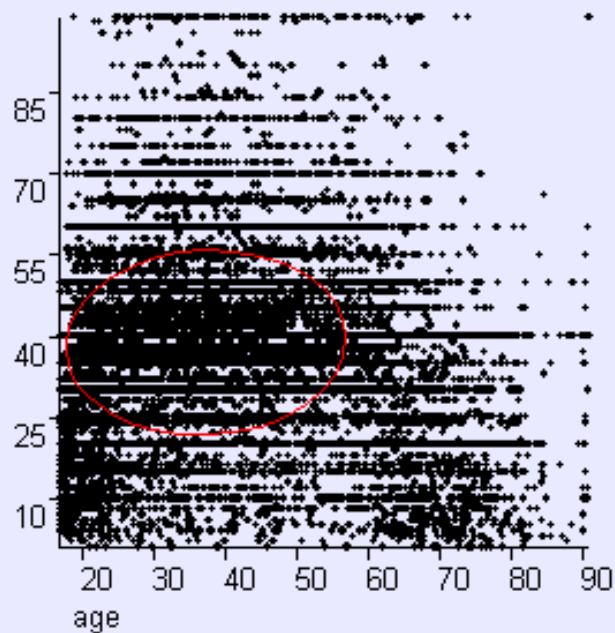
age, hours → wealth

wealth = poor

(prior = 0.760718)

1	mean	cov
age	37.374	198.935 8.70283
hours_worked	38.84	8.70283 152.692

hours_worked

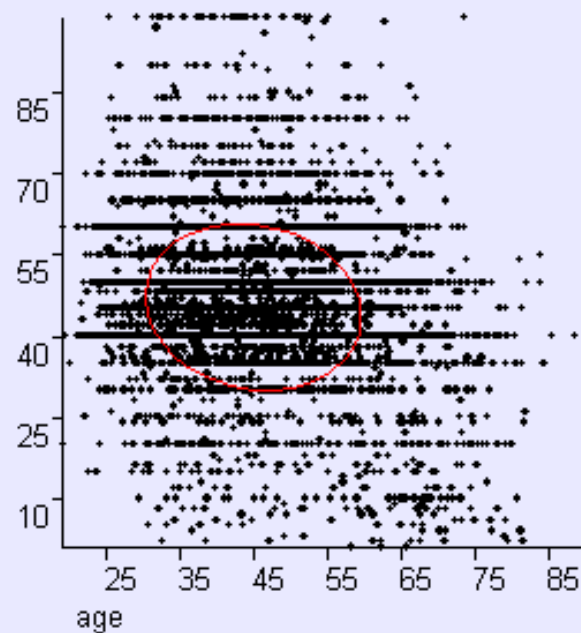


wealth = rich

(prior = 0.239282)

1	mean	cov
age	44.7727	111.618 -14.1565
hours_worked	45.4529	-14.1565 123.014

hours_worked



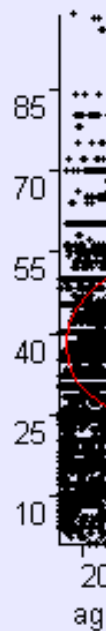
age, hours → wealth

wealth = poor

(prior = 0.760718)

1	mean	cov
age	37.374	198.935 8.70283
hours_worked	38.84	8.70283 152.692

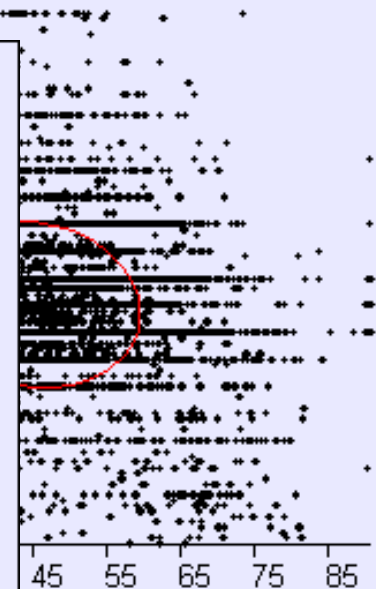
hours_worked



wealth = rich

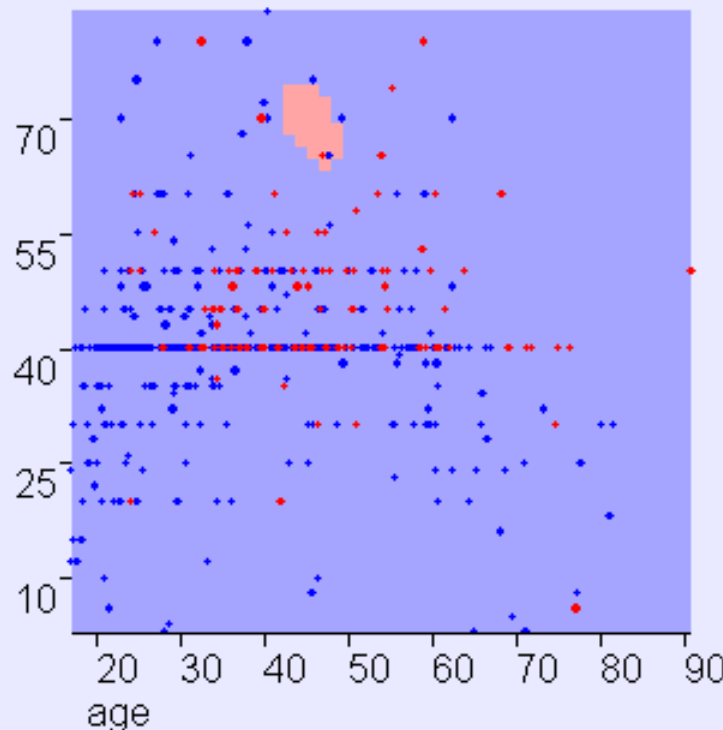
(prior = 0.239282)

1	mean	cov
age	44.7727	111.618 -14.1565
hours_worked	45.4529	-14.1565 123.014



wealth values: poor rich

hours_worked



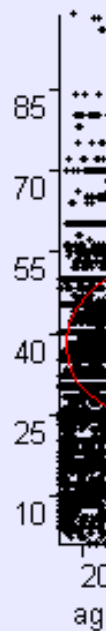
age, hours → wealth

wealth = poor

(prior = 0.760718)

1	mean	cov
age	37.374	198.935 8.70283
hours_worked	38.84	8.70283 152.692

hours_worked



wealth = rich

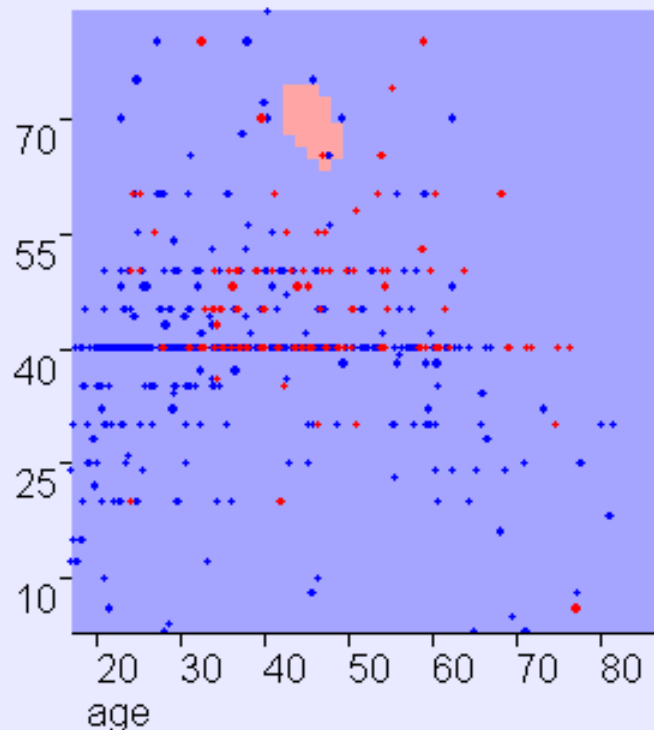
(prior = 0.239282)

1	mean	cov
age	44.7727	111.618 -14.1565
hours_worked	45.4529	-14.1565 123.014



wealth values: poor rich

hours_worked



Having 2 inputs instead of one helps in two ways:

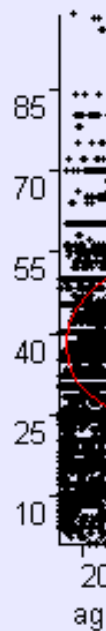
age, hours → wealth

wealth = poor

(prior = 0.760718)

1	mean	cov
age	37.374	198.935 8.70283
hours_worked	38.84	8.70283 152.692

hours_worked



wealth = rich

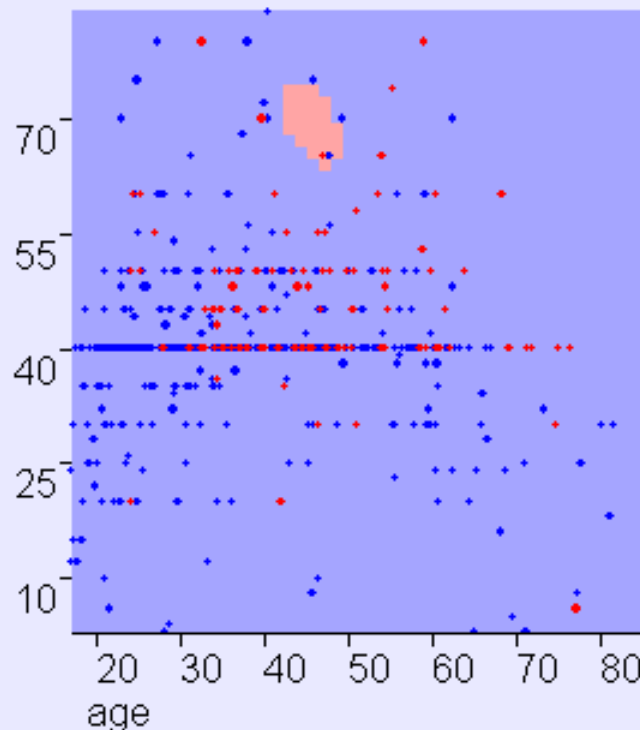
(prior = 0.239282)

1	mean	cov
age	44.7727	111.618 -14.1565
hours_worked	45.4529	-14.1565 123.014



wealth values: poor rich

hours_worked



Having 2 inputs instead of one helps in two ways:

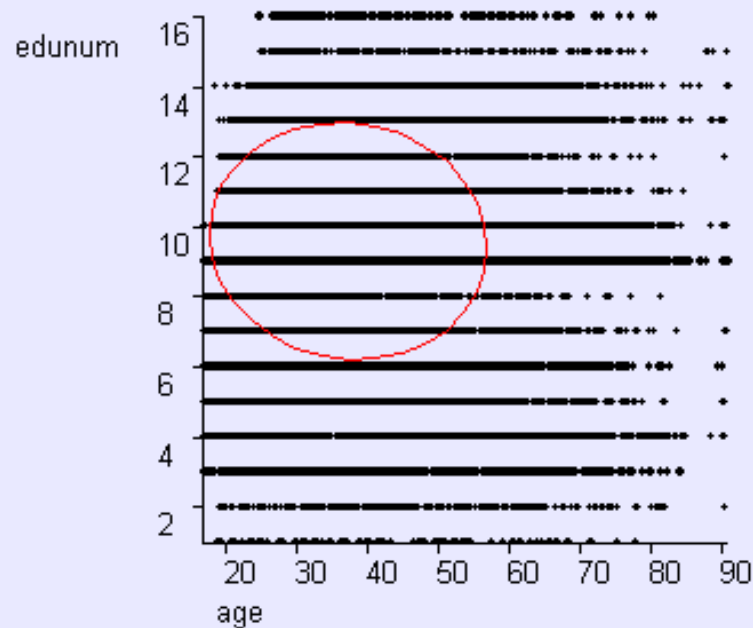
1. Combining evidence from two 1d Gaussians
2. Off-diagonal covariance distinguishes class "shape"

age, edunum \rightarrow wealth

wealth = poor

(prior = 0.760718)

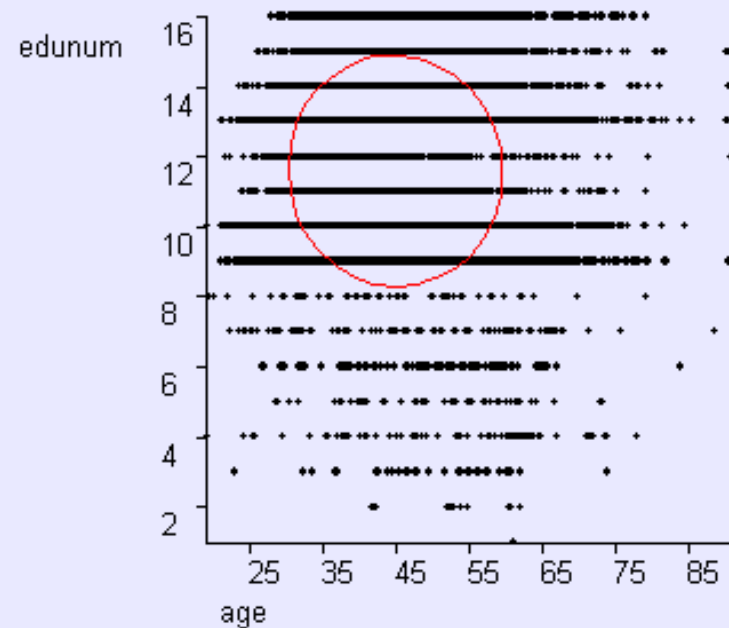
1	mean	cov
age	37.374	198.935 -1.94765
edunum	9.59849	-1.94765 5.94225



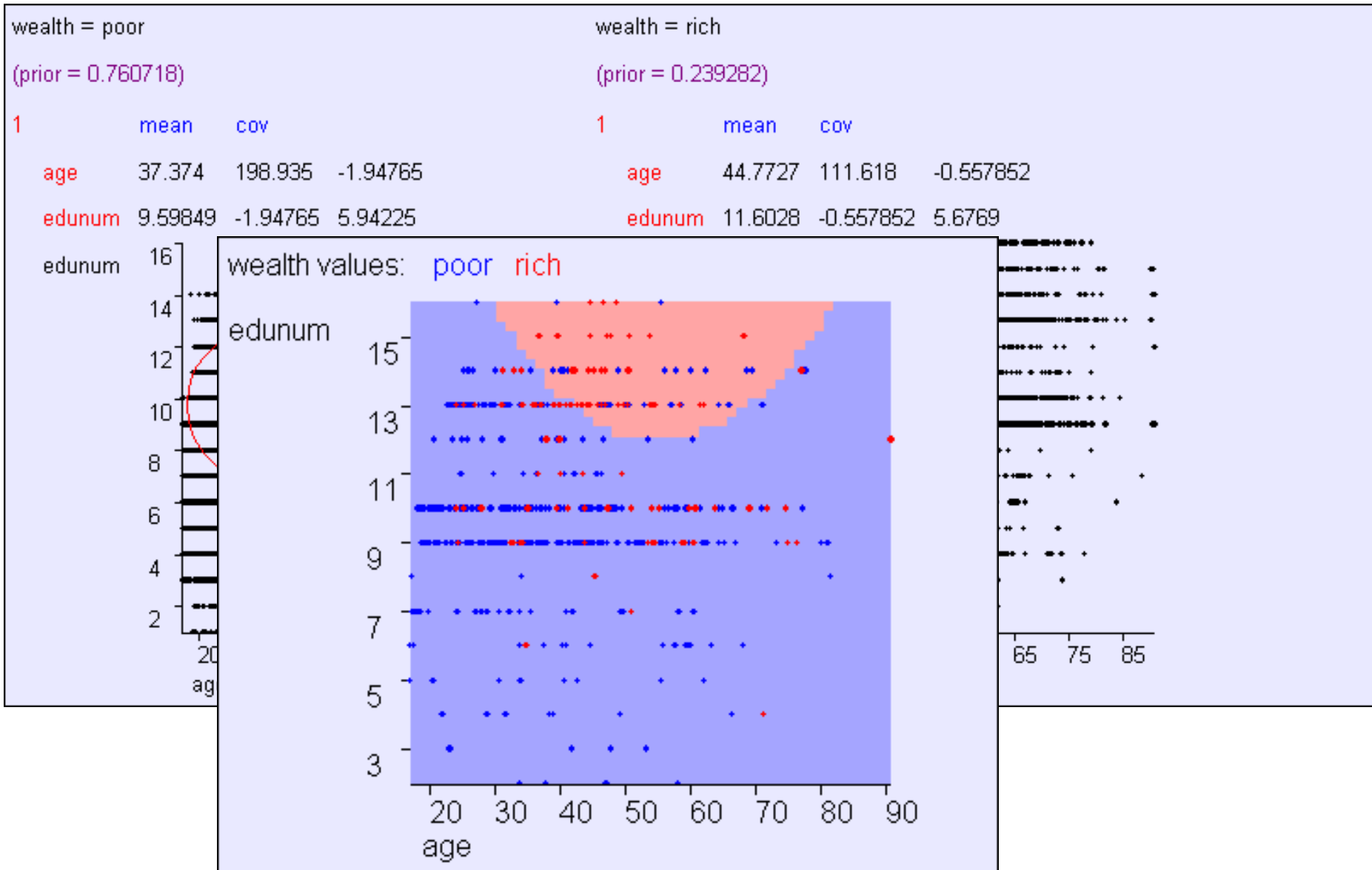
wealth = rich

(prior = 0.239282)

1	mean	cov
age	44.7727	111.618 -0.557852
edunum	11.6028	-0.557852 5.6769



age, edunum → wealth



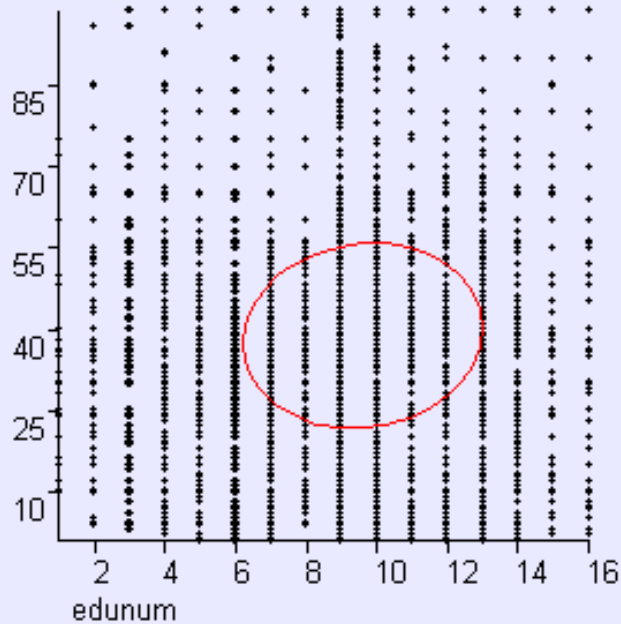
hours, edunum → wealth

wealth = poor

(prior = 0.760718)

1	mean	cov	
edunum	9.59849	5.94225	2.4298
hours_worked	38.84	2.4298	152.692

hours_worked

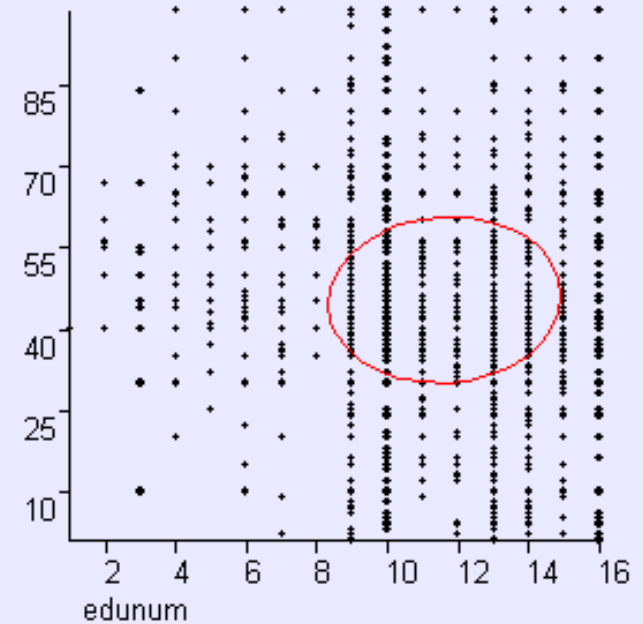


wealth = rich

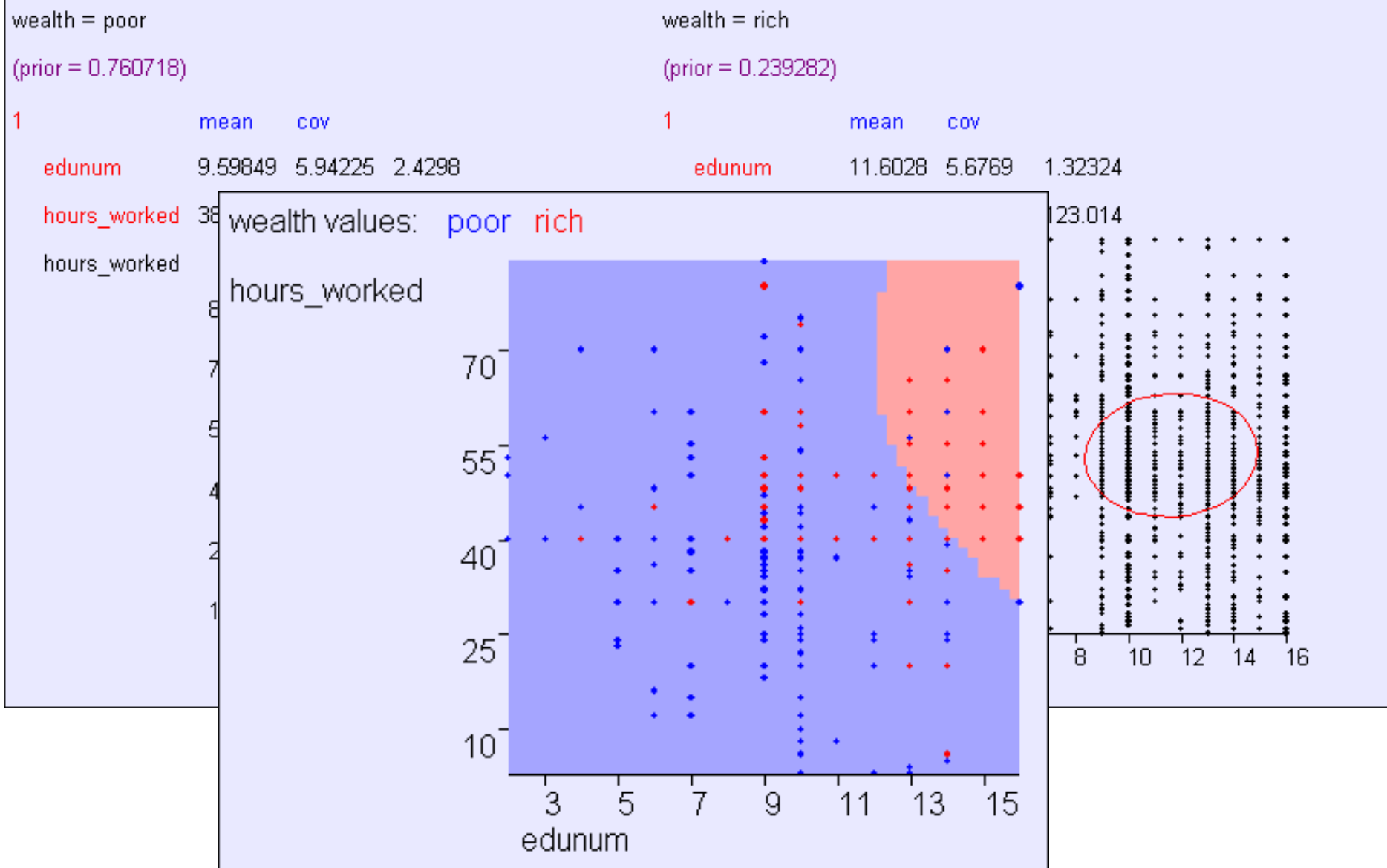
(prior = 0.239282)

1	mean	cov	
edunum	11.6028	5.6769	1.32324
hours_worked	45.4529	1.32324	123.014

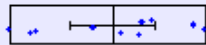
hours_worked

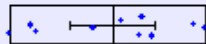


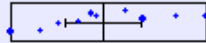
hours, edunum → wealth

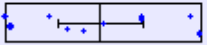


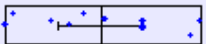
Accuracy

Name	Model	Parameters	FracRight	
age+hours	bayesclass	density=joint submodel=gauss gausstype=general	0.760452 +/- 0.00319521	

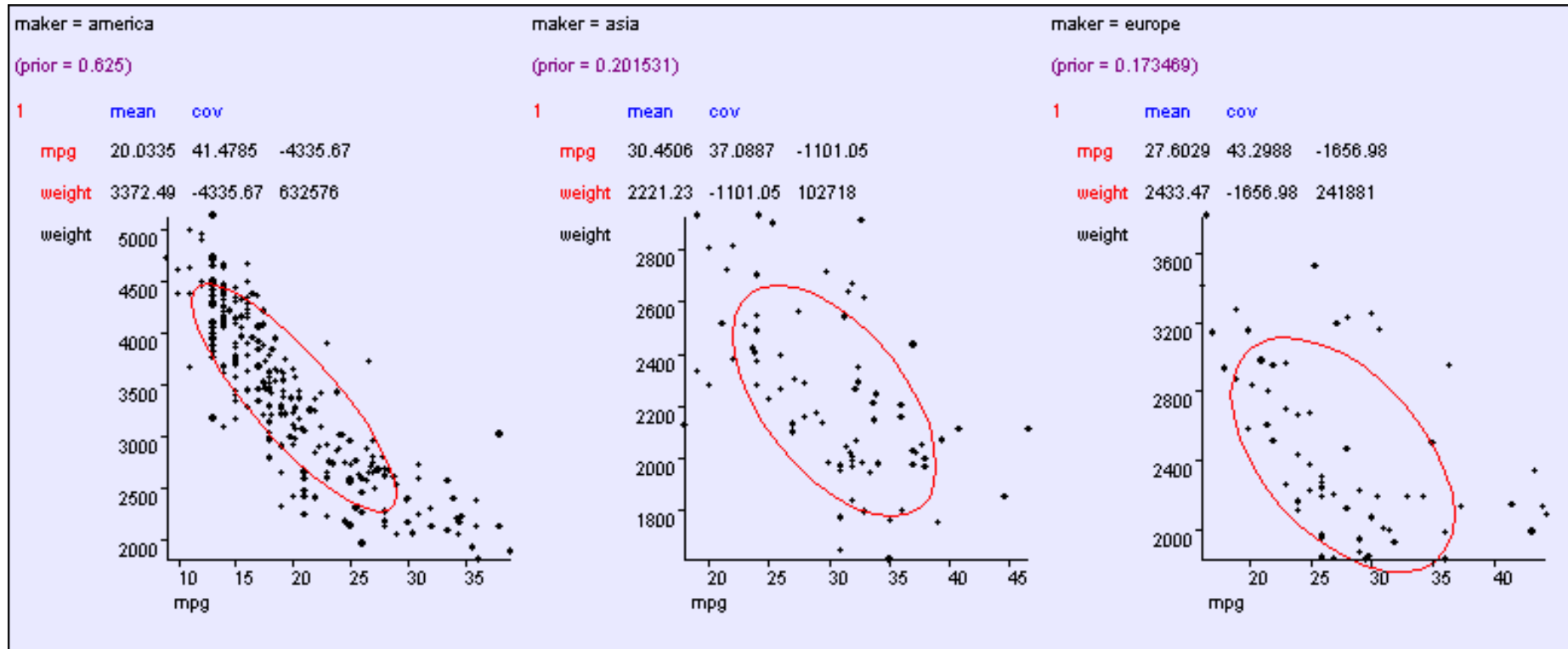
Name	Model	Parameters	FracRight	
age+hours	bayesclass	density=joint submodel=gauss gausstype=general	0.760452 +/- 0.00319521	

Name	Model	Parameters	FracRight	
age+hours+edunum	bayesclass	density=joint submodel=gauss gausstype=general	0.798513 +/- 0.00542432	

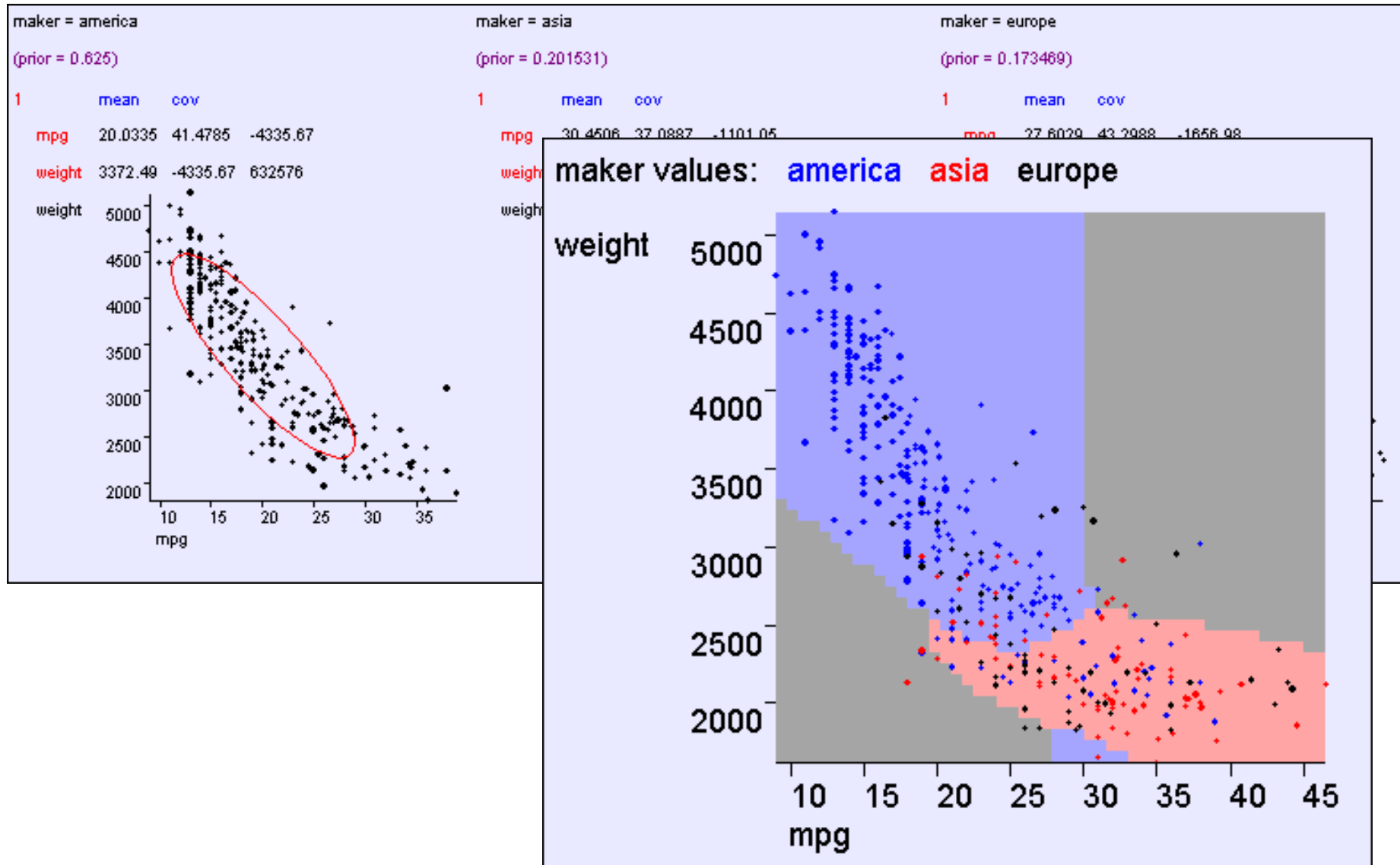
Name	Model	Parameters	FracRight	
a+h+e+capgain	bayesclass	density=joint submodel=gauss gausstype=general	0.793518 +/- 0.00319241	

Name	Model	Parameters	FracRight	
a+h+e+c+taxweight	bayesclass	density=joint submodel=gauss gausstype=general	0.793477 +/- 0.00321524	

An "MPG" example



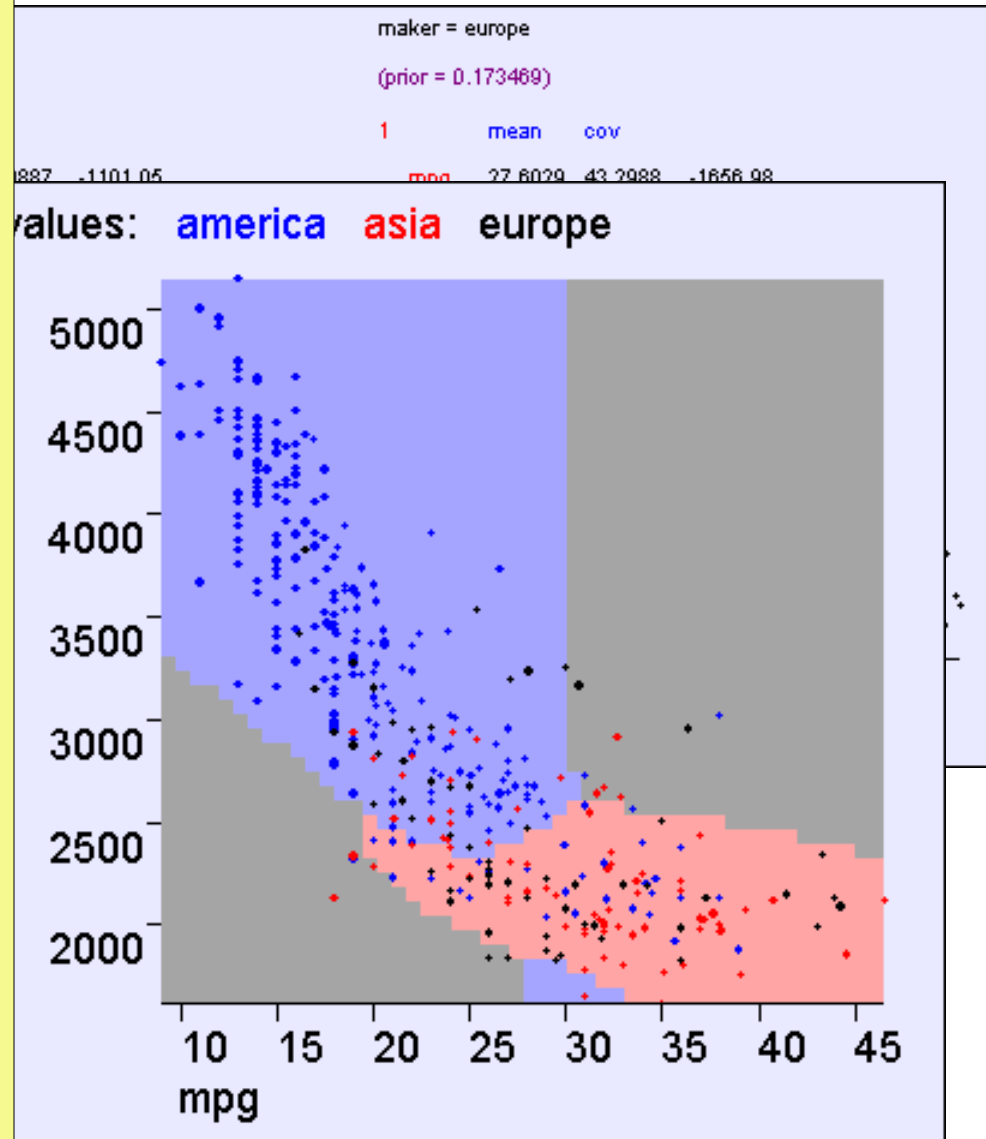
An "MPG" example



An “MPG” example

Things to note:

- Class Boundaries can be weird shapes (hyperconic sections)
- Class regions can be non-simply-connected
- But it's impossible to model arbitrarily weirdly shaped regions
- **Test your understanding:** With one input, must classes be simply connected?



Overfitting dangers

- Problem with “Joint” Bayes classifier:
#parameters exponential with #dimensions.
This means we just memorize the training data, and can overfit.

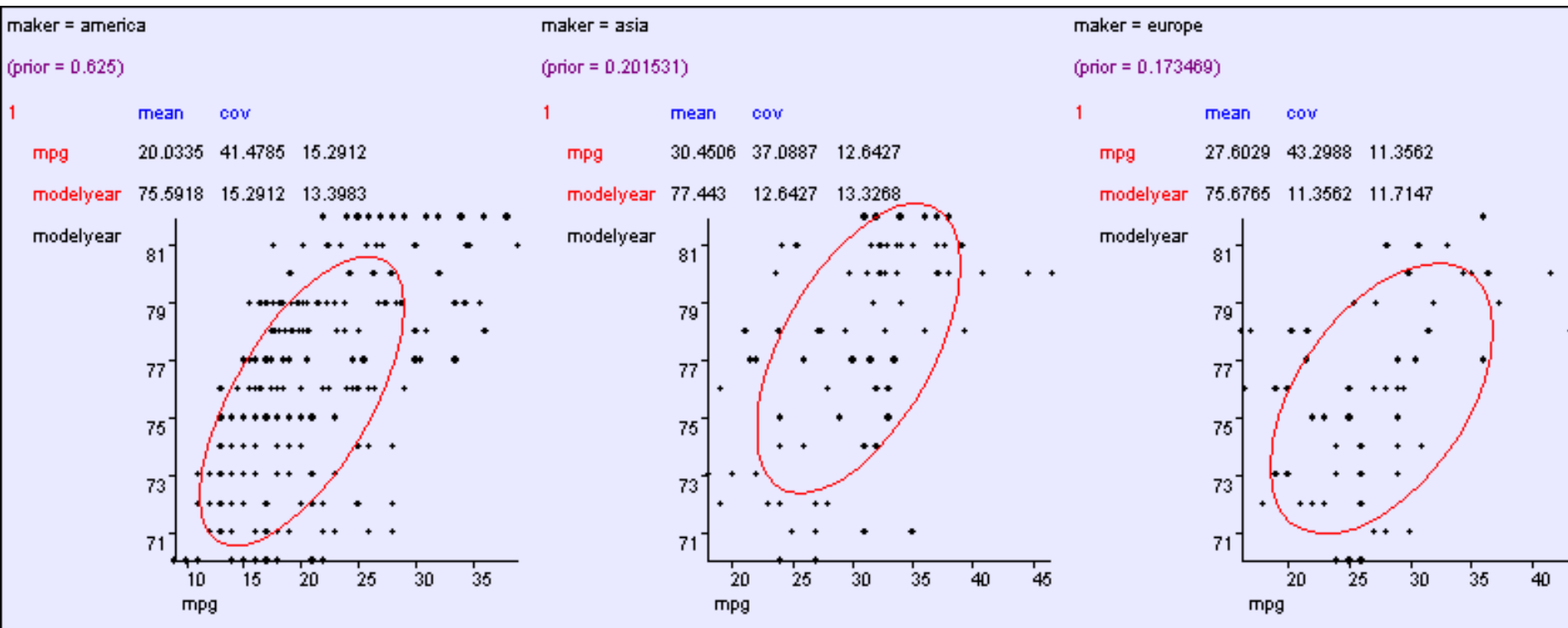
Overfitting dangers

- Problem with “Joint” Bayes classifier:
#parameters exponential with #dimensions.
This means we just memorize the training data, and can overfit.
- Problem with Gaussian Bayes classifier:
#parameters quadratic with #dimensions.
With 10,000 dimensions and only 1,000 datapoints we could overfit.

Question: Any suggested solutions?

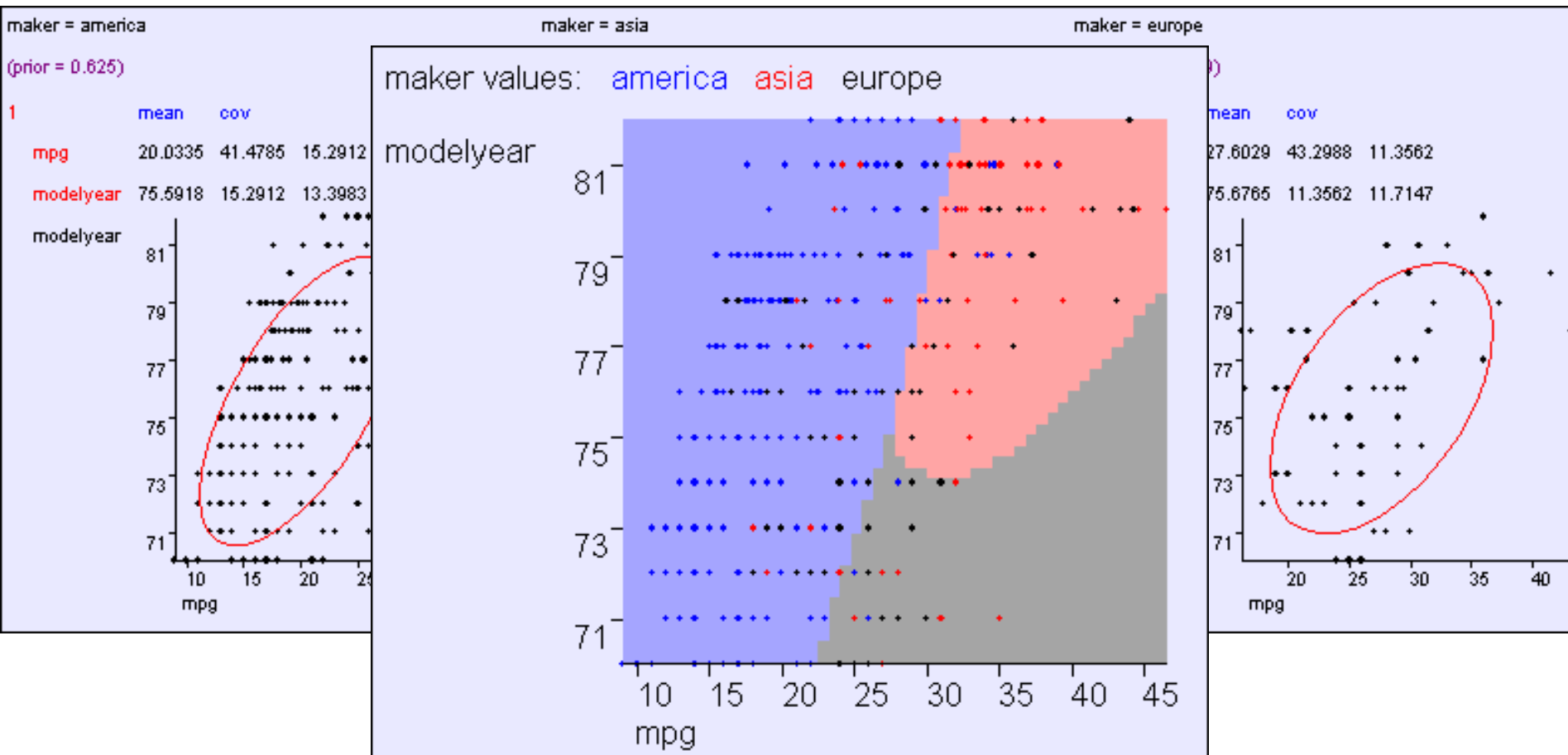
General: $O(m^2)$ parameters

$$\Sigma = \begin{pmatrix} \sigma_{11}^2 & \sigma_{12} & \boxed{?} & \sigma_{1m} \\ \sigma_{12} & \sigma_{22}^2 & \boxed{?} & \sigma_{2m} \\ \boxed{?} & \boxed{?} & \boxed{?} & \boxed{?} \\ \sigma_{1m} & \sigma_{2m} & \boxed{?} & \sigma_{mm}^2 \end{pmatrix}$$



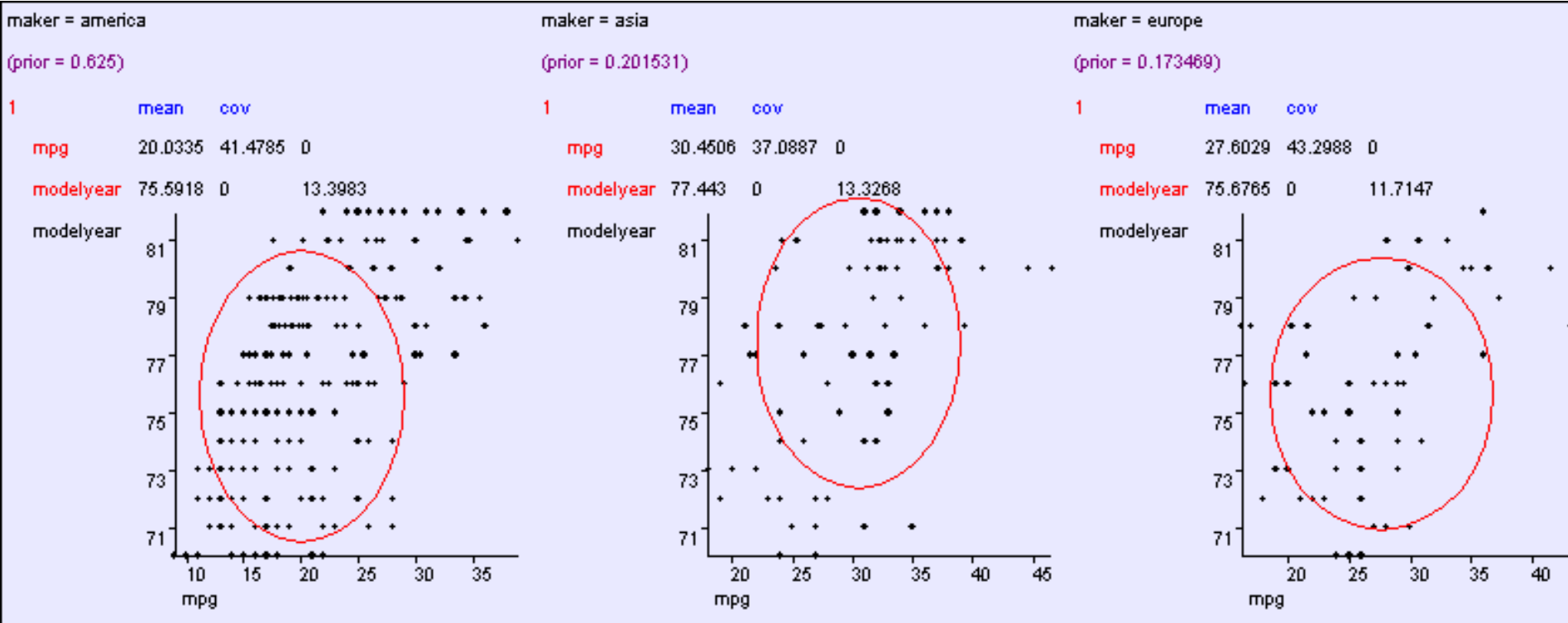
General: $O(m^2)$ parameters

$$\Sigma = \begin{pmatrix} \sigma_{11}^2 & \sigma_{12} & \boxed{?} & \sigma_{1m} \\ \sigma_{12} & \sigma_{22}^2 & \boxed{?} & \sigma_{2m} \\ \boxed{?} & \boxed{?} & \boxed{?} & \boxed{?} \\ \sigma_{1m} & \sigma_{2m} & \boxed{?} & \sigma_m^2 \end{pmatrix}$$



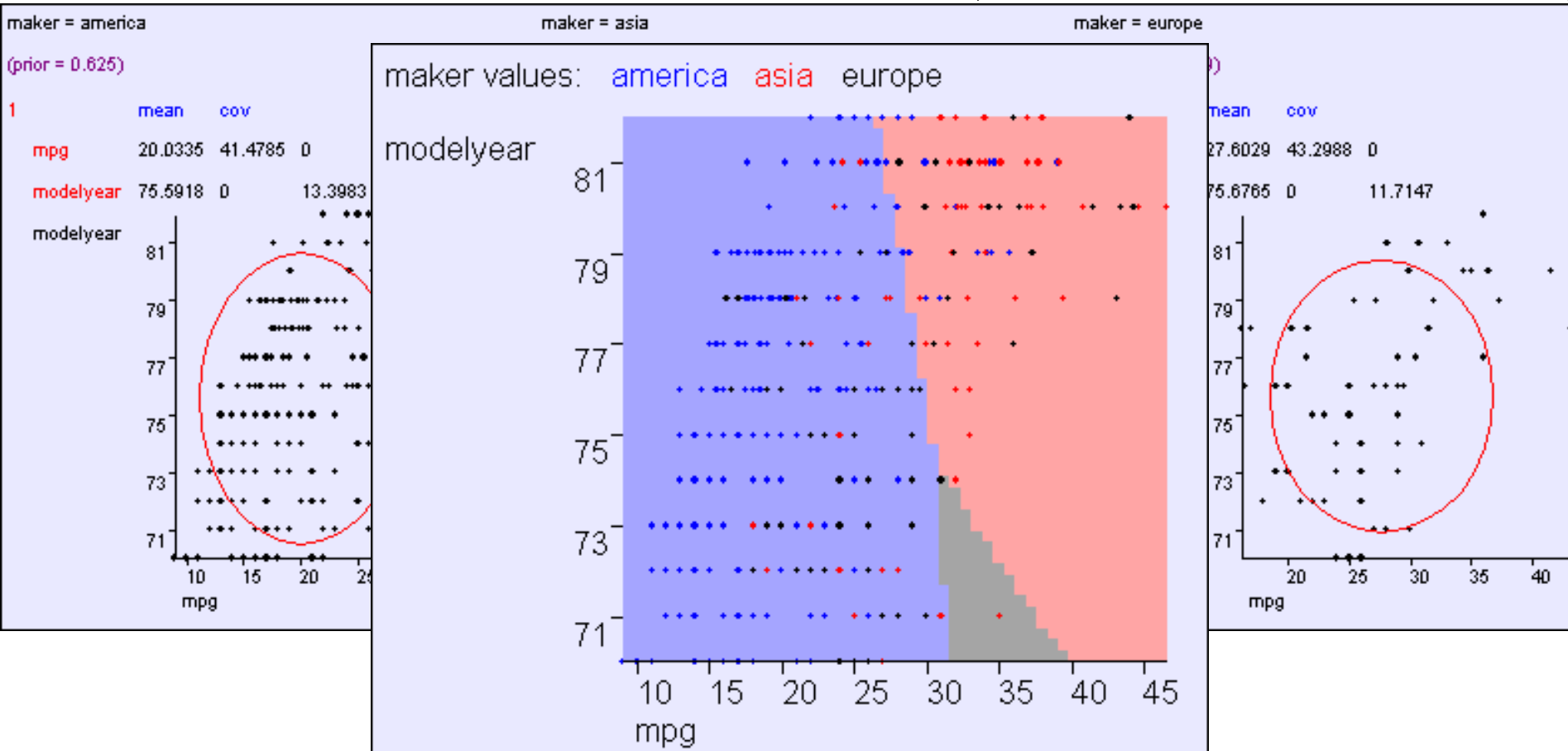
Aligned: $O(m)$ parameters

$$\Sigma = \begin{pmatrix} \sigma^2_1 & 0 & 0 & ? & 0 & 0 \\ 0 & \sigma^2_2 & 0 & ? & 0 & 0 \\ 0 & 0 & \sigma^2_3 & ? & 0 & 0 \\ ? & ? & ? & ? & ? & ? \\ 0 & 0 & 0 & ? & \sigma^2_{m-1} & 0 \\ 0 & 0 & 0 & ? & 0 & \sigma^2_m \end{pmatrix}$$



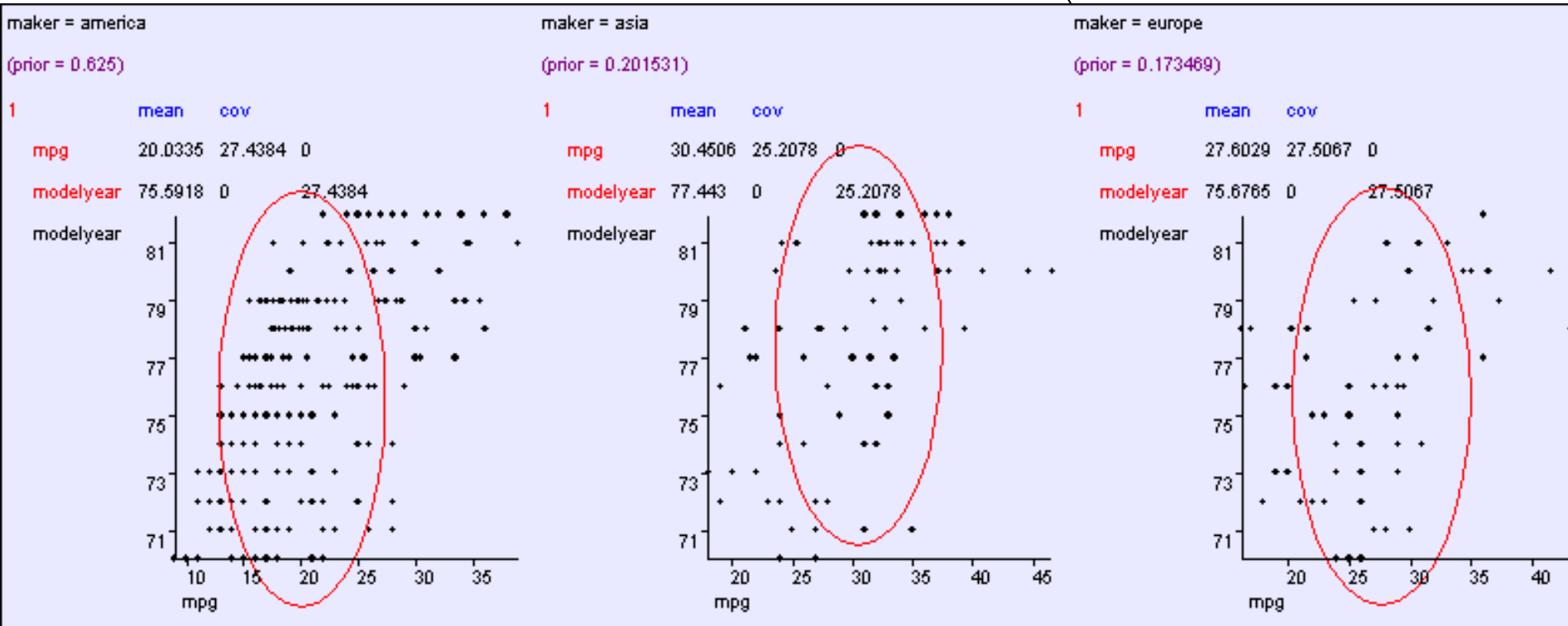
Aligned: $O(m)$ parameters

$$\Sigma = \begin{pmatrix} \sigma^2_1 & 0 & 0 & \boxed{?} & 0 & 0 \\ 0 & \sigma^2_2 & 0 & \boxed{?} & 0 & 0 \\ 0 & 0 & \sigma^2_3 & \boxed{?} & 0 & 0 \\ \boxed{?} & \boxed{?} & \boxed{?} & \boxed{?} & \boxed{?} & \boxed{?} \\ 0 & 0 & 0 & \boxed{?} & \sigma^2_{m-1} & 0 \\ 0 & 0 & 0 & \boxed{?} & 0 & \sigma^2_m \end{pmatrix}$$



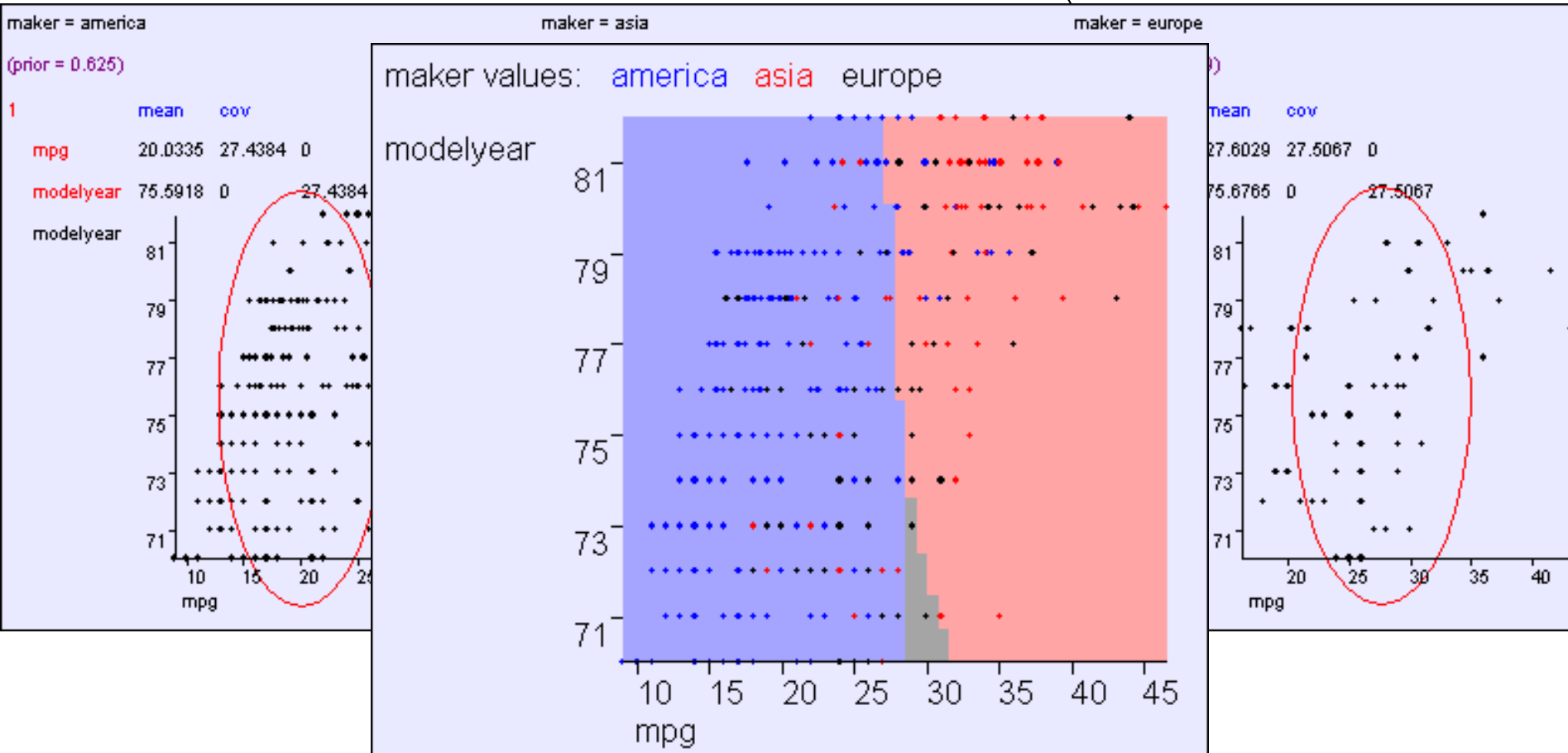
Spherical: O(1) cov parameters

$$\Sigma = \begin{pmatrix} \sigma^2 & 0 & 0 & ? & 0 & 0 \\ 0 & \sigma^2 & 0 & ? & 0 & 0 \\ 0 & 0 & \sigma^2 & ? & 0 & 0 \\ ? & ? & ? & ? & ? & ? \\ 0 & 0 & 0 & ? & \sigma^2 & 0 \\ 0 & 0 & 0 & ? & 0 & \sigma^2 \end{pmatrix}$$



Spherical: O(1) cov parameters

$$\Sigma = \begin{pmatrix} \sigma^2 & 0 & 0 & \boxed{?} & 0 & 0 \\ 0 & \sigma^2 & 0 & \boxed{?} & 0 & 0 \\ 0 & 0 & \sigma^2 & \boxed{?} & 0 & 0 \\ \boxed{?} & \boxed{?} & \boxed{?} & \boxed{?} & \boxed{?} & \boxed{?} \\ 0 & 0 & 0 & \boxed{?} & \sigma^2 & 0 \\ 0 & 0 & 0 & \boxed{?} & 0 & \sigma^2 \end{pmatrix}$$



BCs that have both real and categorical inputs?

		Categorical inputs only	Real-valued inputs only	Mixed Real / Cat okay
<div>Inputs</div> <div><div>Classifier</div><div>Predict category</div></div>	Joint BC Naïve BC	Gauss BC	Dec Tree BC Here???	
<div>Inputs</div> <div><div>Density Estimator</div><div>Prob-ability</div></div>	Joint DE Naïve DE	Gauss DE		
<div>Inputs</div> <div><div>Regressor</div><div>Predict real no.</div></div>				

BCs that have both real and categorical inputs?

		Categorical inputs only	Real-valued inputs only	Mixed Real / Cat okay
<div>Inputs</div> <div><div>Classifier</div><div>Predict category</div></div>	Joint BC Naïve BC	Gauss BC	Dec Tree BC Here???	
<div>Inputs</div> <div><div>Density Estimator</div><div>Prob-ability</div></div>	Joint DE Naïve DE	Gauss DE	<div>Easy!</div> <div>Guess how?</div>	
<div>Inputs</div> <div><div>Regressor</div><div>Predict real no.</div></div>				

BCs that have both real and categorical inputs?

		Categorical inputs only	Real-valued inputs only	Mixed Real / Cat okay
Inputs	<div>Classifier</div> <div>Predict category</div>	Joint BC Naïve BC	Gauss BC	Dec Tree Gauss/Joint BC Gauss Naïve BC
Inputs	<div>Density Estimator</div> <div>Prob-ability</div>	Joint DE Naïve DE	Gauss DE Gauss DE	Gauss/Joint DE Gauss Naïve DE
Inputs	<div>Regressor</div> <div>Predict real no.</div>			

BCs that have both real and categorical inputs?

		Categorical inputs only	Real-valued inputs only	Mixed Real / Cat okay
Inputs	<div>Classifier</div> <div>Predict category</div>	Joint BC Naïve BC	Gauss BC	Dec Tree Gauss/Joint BC Gauss Naïve BC
Inputs	<div>Density Estimator</div> <div>Prob-ability</div>	Joint DE Naïve DE	Gauss DE Gauss DE	Gauss/Joint DE Gauss Naïve DE
Inputs	<div>Regressor</div> <div>Predict real no.</div>			

Mixed Categorical / Real Density Estimation

- Write $\mathbf{x} = (\mathbf{u}, \mathbf{v}) = (\underbrace{u_1, u_2, \dots, u_q}_{\text{Real valued}}, \underbrace{v_1, v_2, \dots, v_{m-q}}_{\text{Categorical valued}})$

$$P(\mathbf{x} | M) = P(\mathbf{u}, \mathbf{v} | M)$$

(where **M** is any Density Estimation **M**odel)

Not sure which tasty
DE to enjoy? Try our...

Joint / Gauss DE Combo

$$P(\mathbf{u}, \mathbf{v} \mid M) = P(\mathbf{u} \mid \mathbf{v}, M) P(\mathbf{v} \mid M)$$

Gaussian with
parameters
depending on \mathbf{v}

Big “m-q”-dimensional
lookup table

MLE learning of the Joint / Gauss DE Combo

$$P(\mathbf{u}, \mathbf{v} \mid M) = P(\mathbf{u} \mid \mathbf{v}, M) P(\mathbf{v} \mid M)$$

μ_v = Mean of \mathbf{u} among
records matching \mathbf{v}

Σ_v = Cov. of \mathbf{u} among
records matching \mathbf{v}

q_v = Fraction of records
that match \mathbf{v}

$$\mathbf{u} \mid \mathbf{v}, M \sim N(\mu_v, \Sigma_v) \quad , \quad P(\mathbf{v} \mid M) = q_v$$

MLE learning of the Joint / Gauss DE Combo

$$P(\mathbf{u}, \mathbf{v} \mid M) = P(\mathbf{u} \mid \mathbf{v}, M) P(\mathbf{v} \mid M)$$

$$\mu_{\mathbf{v}} = \text{Mean of } \mathbf{u} \text{ among records matching } \mathbf{v}$$

$$\Sigma_{\mathbf{v}} = \text{Cov. of } \mathbf{u} \text{ among records matching } \mathbf{v}$$

$$q_{\mathbf{v}} = \text{Fraction of records that match } \mathbf{v}$$

$$\frac{1}{R_{\mathbf{v}}} \sum_{k \text{ s.t. } \mathbf{v}_k = \mathbf{v}} \mathbf{u}_k$$

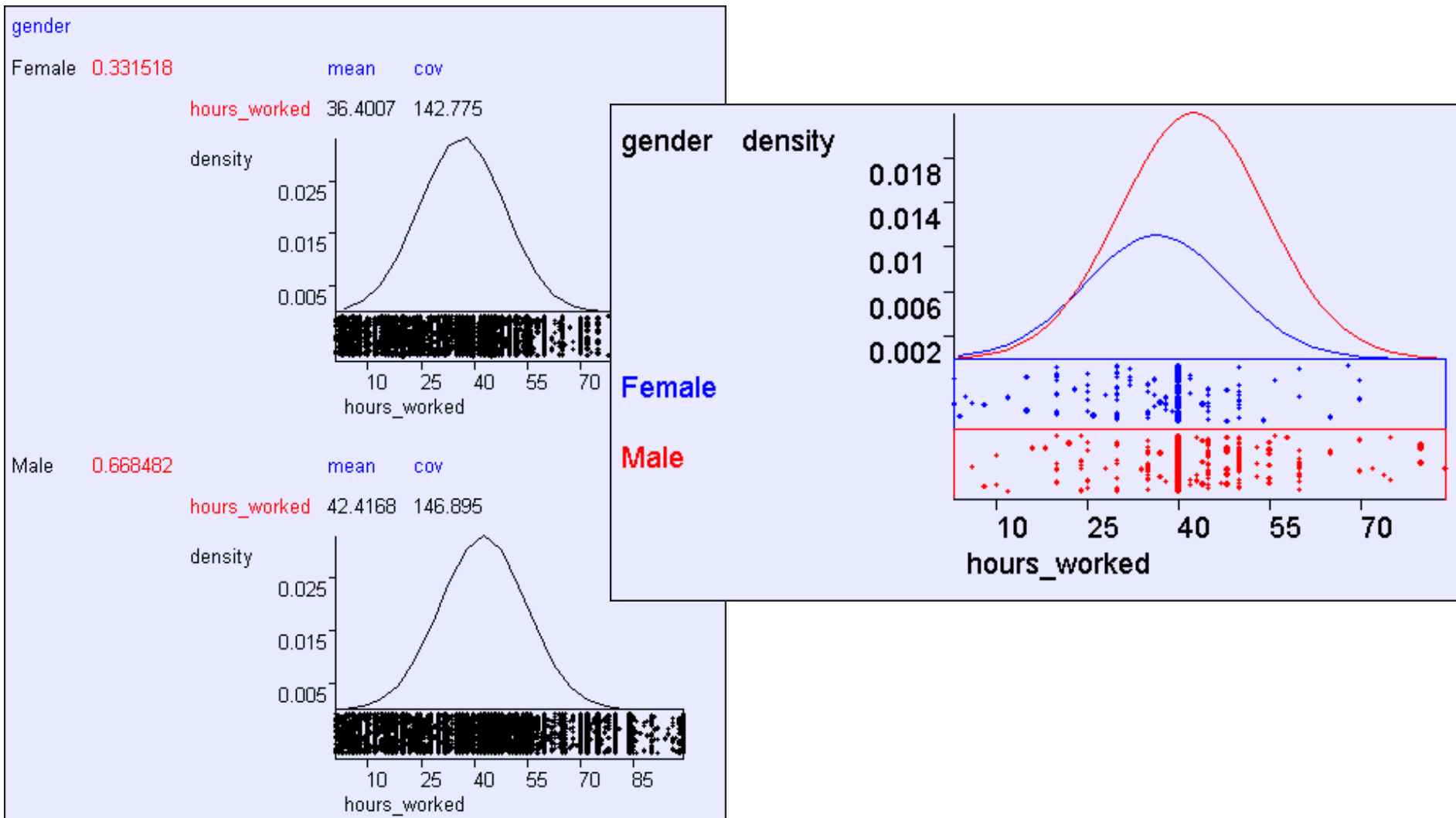
$$\frac{1}{R_{\mathbf{v}}} \sum_{k \text{ s.t. } \mathbf{v}_k = \mathbf{v}} (\mathbf{u}_k - \mu_{\mathbf{v}})(\mathbf{u}_k - \mu_{\mathbf{v}})^T$$

$$\frac{R_{\mathbf{v}}}{R}$$


$$R_{\mathbf{v}} = \# \text{ records that match } \mathbf{v}$$


$$\mathbf{u} \mid \mathbf{v}, M \sim N(\mu_{\mathbf{v}}, \Sigma_{\mathbf{v}}), \quad P(\mathbf{v} \mid M) = q_{\mathbf{v}}$$

Gender and Hours Worked*



*As with all the results from the UCI "adult census" dataset, we can't draw any real-world conclusions since it's such a non-real-world sample

What we just did  Joint / Gauss DE
Combo

What we do next  Joint / Gauss **BC**
Combo

Joint / Gauss **BC** Combo

$$\begin{aligned} P(Y = i \mid \mathbf{u}, \mathbf{v}) &= \frac{p(\mathbf{u}, \mathbf{v} \mid M_i) P(Y = i)}{p(\mathbf{u}, \mathbf{v})} \\ &= \frac{p(\mathbf{u}, \mid \mathbf{v}, M_i) p(\mathbf{v} \mid M_i) P(Y = i)}{p(\mathbf{u}, \mathbf{v})} \\ &= \frac{N(\mathbf{u}; \boldsymbol{\mu}_{i,\mathbf{v}}, \boldsymbol{\Sigma}_{i,\mathbf{v}}) q_{i,\mathbf{v}} p_i}{p(\mathbf{u}, \mathbf{v})} \end{aligned}$$

Joint / Gauss **BC** Combo

$$P(Y = i | \mathbf{u}, \mathbf{v}) = \frac{p(\mathbf{u}, \mathbf{v} | M_i) P(Y = i)}{p(\mathbf{u}, \mathbf{v})}$$

$\mu_{i,v}$	= Mean of \mathbf{u} among records matching \mathbf{v} and in which $y=i$
-------------	---

$$= \frac{p(\mathbf{u}, | \mathbf{v}, M_i) p(\mathbf{v} | M_i) P(Y = i)}{p(\mathbf{u}, \mathbf{v})}$$

$\Sigma_{i,v}$	= Cov. of \mathbf{u} among records matching \mathbf{v} and in which $y=i$
----------------	---

$$= \frac{N(\mathbf{u}; \mu_{i,v}, \Sigma_{i,v}) q_{i,v} p_i}{p(\mathbf{u}, \mathbf{v})}$$

$q_{i,v}$	= Fraction of “ $y=i$ ” records that match \mathbf{v}
-----------	---

p_i	= Fraction of records that match “ $y=i$ ”
-------	--

Rather so-so-notation for
“Gaussian with mean $\mu_{i,v}$ and
covariance $\Sigma_{i,v}$ evaluated at \mathbf{u} ”

Gender, Hours→Wealth

wealth = poor

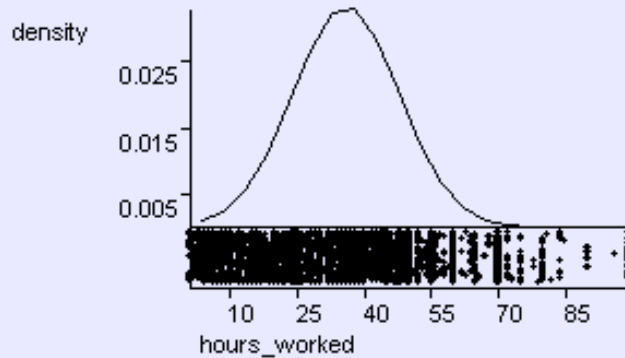
(prior = 0.760718)

gender

Female 0.388185

mean cov

hours_worked 35.876 140.542



wealth = rich

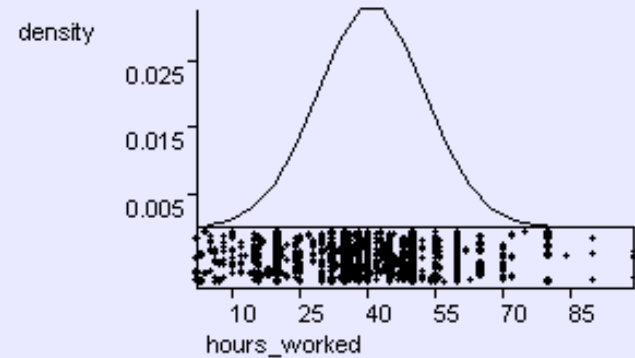
(prior = 0.239282)

gender

Female 0.151365

mean cov

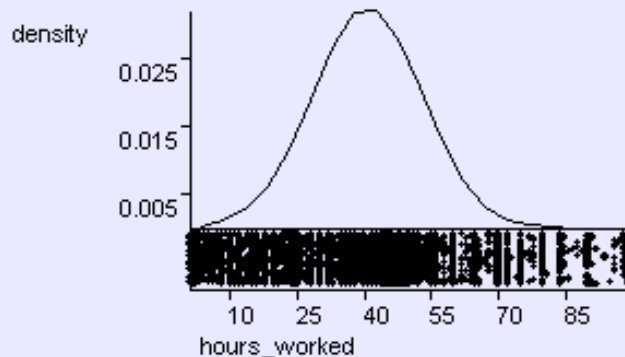
hours_worked 40.6789 140.511



Male 0.611815

mean cov

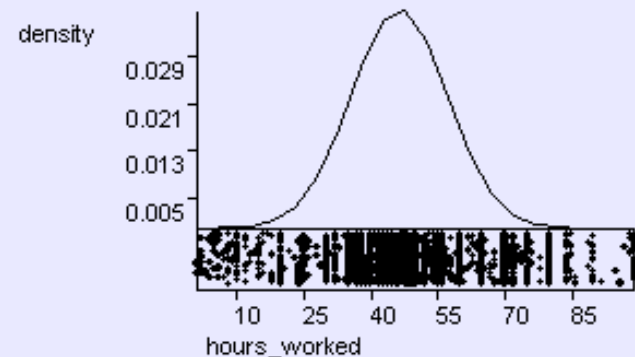
hours_worked 40.7207 151.295



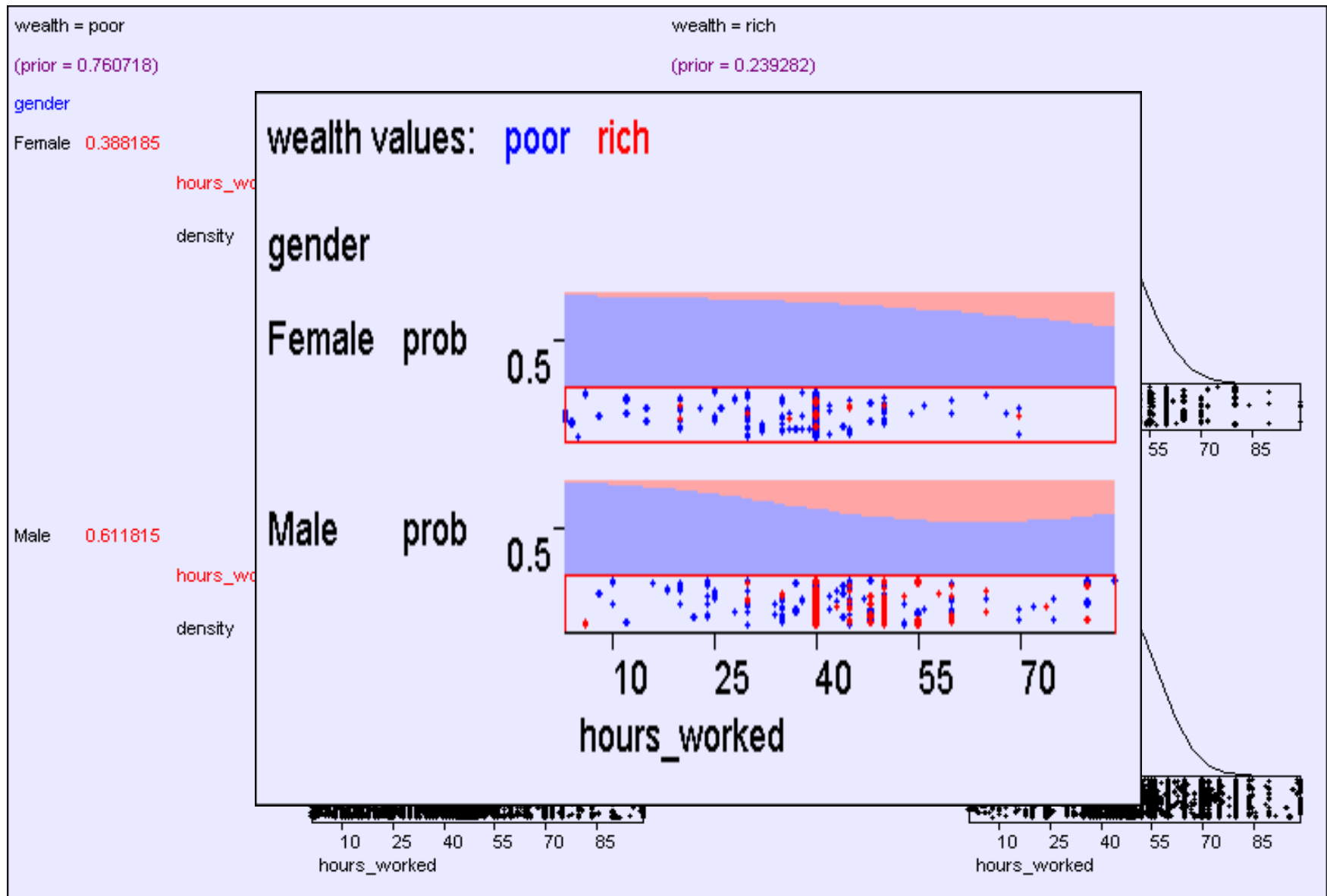
Male 0.848635

mean cov

hours_worked 46.3044 115.117



Gender, Hours→Wealth

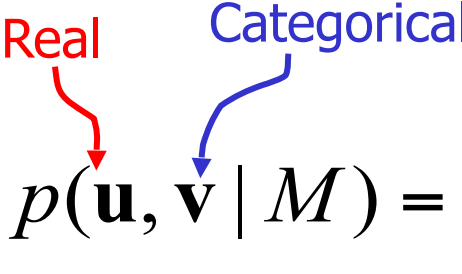


Joint / Gauss DE Combo and Joint / Gauss BC Combo: The downside

- (Yawn...we've done this before...)
More than a few categorical attributes blah blah
blah massive table blah blah lots of parameters
blah blah just memorize training data blah blah
blah do worse on future data blah blah need to
be more conservative blah

Naïve/Gauss combo for Density Estimation

Real Categorical


$$p(\mathbf{u}, \mathbf{v} \mid M) = \left(\prod_{j=1}^q p(u_j \mid M) \right) \left(\prod_{j=1}^{m-q} P(v_j \mid M) \right)$$
$$u_j \mid M \sim N(\mu_j, \sigma_j^2) \quad v_j \mid M \sim \text{Multinomial}[q_{j1}, q_{j2}, \dots, q_{jN_j}]$$

How many parameters?

Naïve/Gauss combo for Density Estimation

Real Categorical

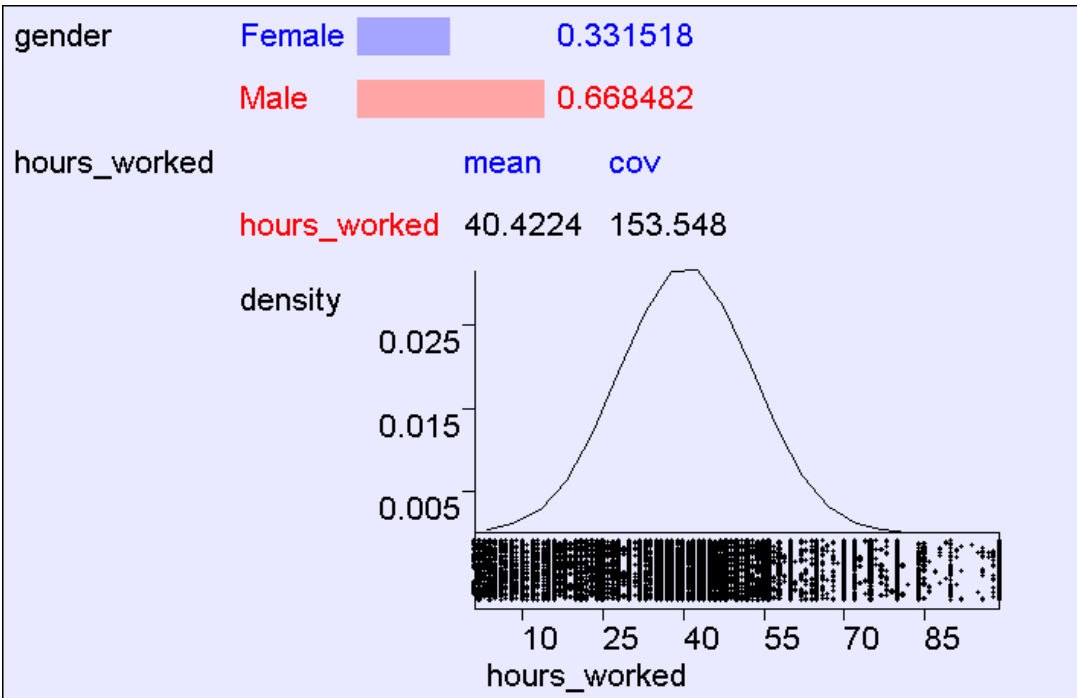
$$p(\mathbf{u}, \mathbf{v} \mid M) = \left(\prod_{j=1}^q p(u_j \mid M) \right) \left(\prod_{j=1}^{m-q} P(v_j \mid M) \right)$$
$$u_j \mid M \sim N(\mu_j, \sigma_j^2) \quad v_j \mid M \sim \text{Multinomial}[q_{j1}, q_{j2}, \dots, q_{jN_j}]$$

$$\mu_j = \frac{1}{R} \sum_k u_{kj}$$

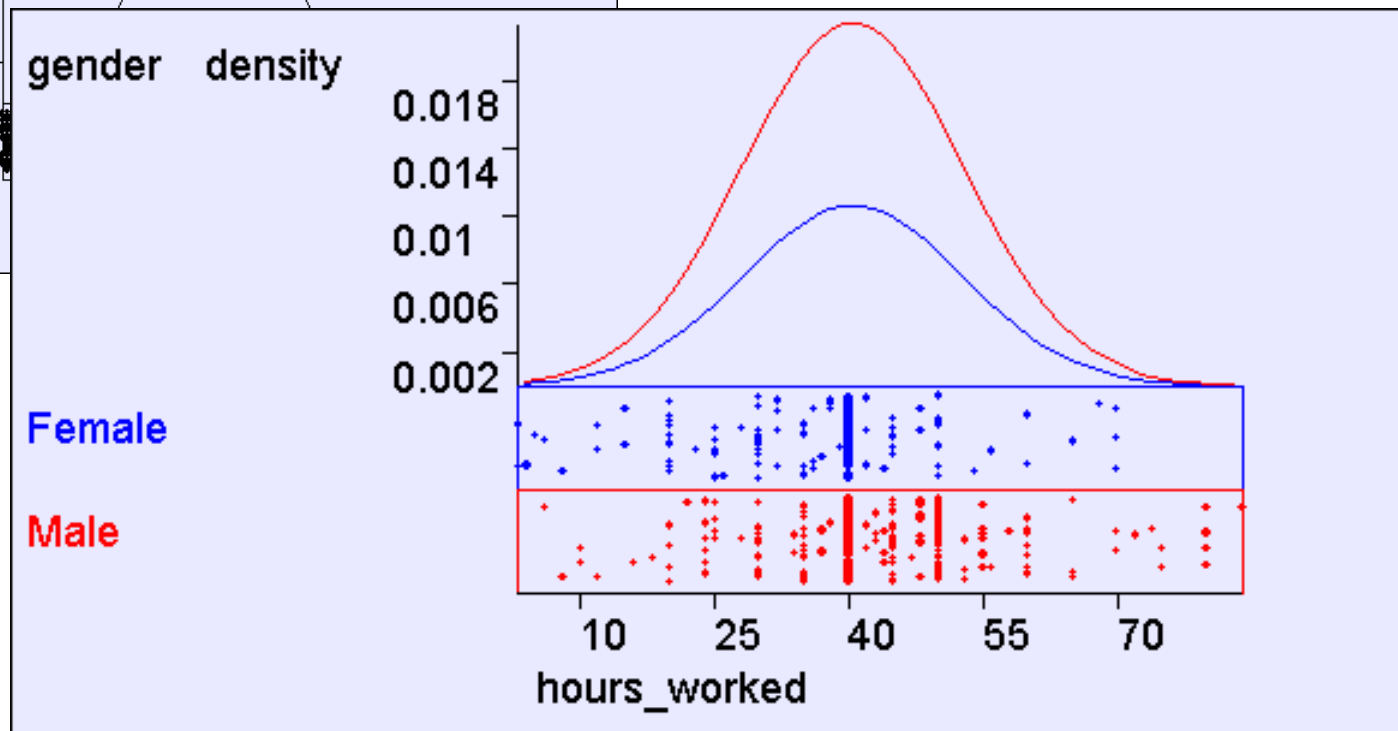
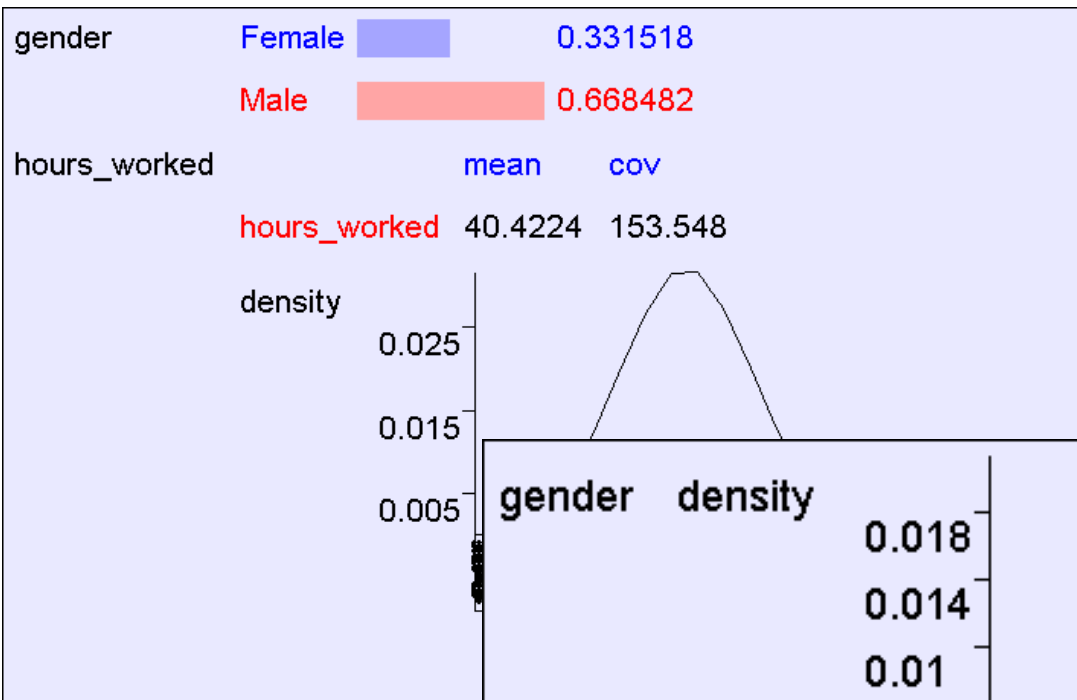
$$\sigma_j^2 = \frac{1}{R} \sum_k (u_{kj} - \mu_j)^2$$

$$q_{jh} = \frac{\# \text{ of records in which } v_j = h}{R}$$

Naïve/Gauss DE Example



Naïve/Gauss DE Example



Naïve / Gauss BC

$$P(Y = i | \mathbf{u}, \mathbf{v}) = \frac{p(\mathbf{u}, \mathbf{v} | Y = i)P(Y = i)}{p(\mathbf{u}, \mathbf{v})}$$

$$= \frac{1}{p(\mathbf{u}, \mathbf{v})} \prod_{j=1}^q p(u_j | \mu_{ij}, \sigma_{ij}^2) \prod_{j=1}^{m-q} P(v_j | \mathbf{q}_{ij}) P(Y = i)$$

$$= \frac{1}{p(\mathbf{u}, \mathbf{v})} \prod_{j=1}^q N(u_j; \mu_{ij}, \sigma_{ij}^2) \prod_{j=1}^{m-q} q_{ij}[v_j] p_i$$

μ_{ij} = Mean of u_j among records in which $y=i$

σ_{ij}^2 = Var. of u_j among records in which $y=i$

$q_{ij}[h]$ = Fraction of “ $y=i$ ” records in which $v_j = h$

p_i = Fraction of records that match “ $y=i$ ”

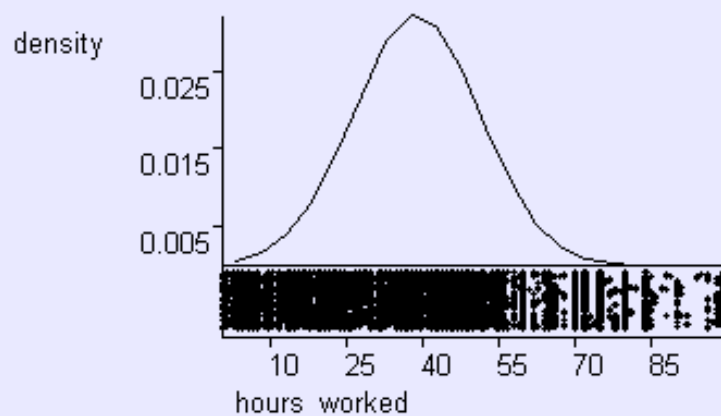
Gauss / Naïve BC Example

wealth = poor

(prior = 0.760718)

gender Female 0.388185
 Male 0.611815

hours_worked mean cov
hours_worked 38.84 152.692

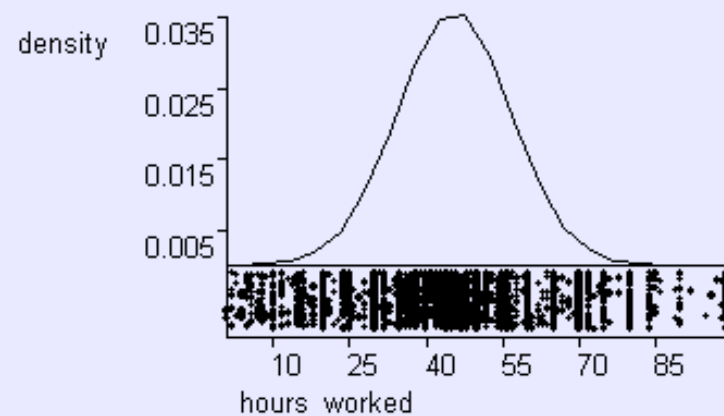


wealth = rich

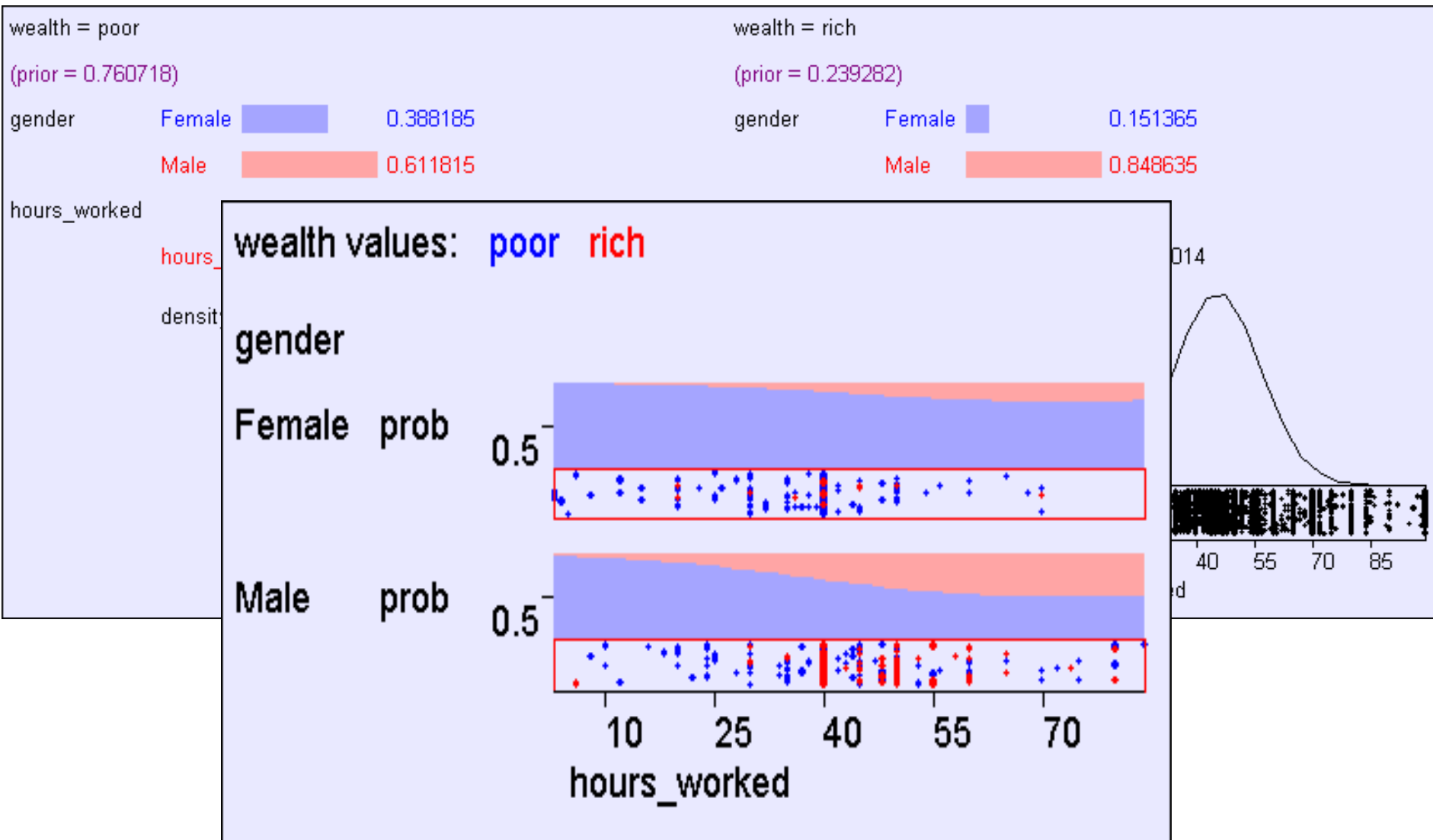
(prior = 0.239282)

gender Female 0.151365
 Male 0.848635


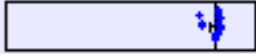

hours_worked mean cov
hours_worked 45.4529 123.014






Gauss / Naïve BC Example






Learn Wealth from 15 attributes

Name	Model	Parameters	FracRight	
Model1	bayesclass	density=joint submodel=gauss gausstype=general	0.718009 +/- 0.00570714	
Model2	bayesclass	density=naive submodel=gauss gausstype=general	0.832234 +/- 0.00288377	
Model3	dtree	max_children=4 ne_splits=y max_pchance=0.05 adjust_chi=y max_nodes=50	0.850702 +/- 0.00364538	




Learn Wealth from 15 attributes

Name	Model	Parameters	FracRight	
Model1	bayesclass	density=joint submodel=gauss gausstype=general	0.718009 +/- 0.00570714	
Model2	bayesclass	density=naive submodel=gauss gausstype=general	0.832234 +/- 0.00288377	
Model3	dtree	max_children=4 ne_splits=y max_pchance=0.05 adjust_chi=y max_nodes=50	0.850702 +/- 0.00364538	

Same data, except all
real values discretized
to 3 levels

Model	Parameters	FracRight	
bayesclass	density=joint submodel=gauss gausstype=general	0.800418 +/- 0.00321903	
bayesclass	density=naive submodel=gauss gausstype=general	0.819745 +/- 0.00240386	
dtree	max_children=4 ne_splits=y max_pchance=0.05 adjust_chi=y max_nodes=50	0.826113 +/- 0.00327583	

Learn Race from 15 attributes

Name	Model	Parameters	FracRight	
Model1	bayesclass	density=joint submodel=gauss gausstype=general	0.391303 +/- 0.00586792	
Model2	bayesclass	density=naive submodel=gauss gausstype=general	0.788686 +/- 0.00560675	
Model3	dtree	max_children=4 ne_splits=y max_pchance=0.05 adjust_chi=y max_nodes=50	0.860919 +/- 0.00272011	

What you should know

- A lot of this should have just been a corollary of what you already knew
- Turning Gaussian DEs into Gaussian BCs
- Mixing Categorical and Real-Valued

Questions to Ponder

- Suppose you wanted to create an example dataset where a BC involving Gaussians crushed decision trees like a bug. What would you do?
- Could you combine Decision Trees and Bayes Classifiers? How? (maybe there is more than one possible way)