

Perspectives of Predictive Modeling

Arthur Charpentier

charpentier.arthur@uqam.ca

[http ://freakonometrics.hypotheses.org/](http://freakonometrics.hypotheses.org/)



(SOA Webcast, November 2013)

Agenda

- **Introduction to Predictive Modeling**
 - Prediction, best estimate, expected value and confidence interval
 - Parametric versus nonparametric models
- **Linear Models and (Ordinary) Least Squares**
 - From least squares to the Gaussian model
 - Smoothing continuous covariates
- **From Linear Models to G.L.M.**
- **Modeling a TRUE-FALSE variable**
 - The logistic regression
 - R.O.C. curve
 - Classification tree (and random forests)
- **From individual to functional data**

Prediction ? Best estimate ?

E.g. predicting someone's weight (Y)

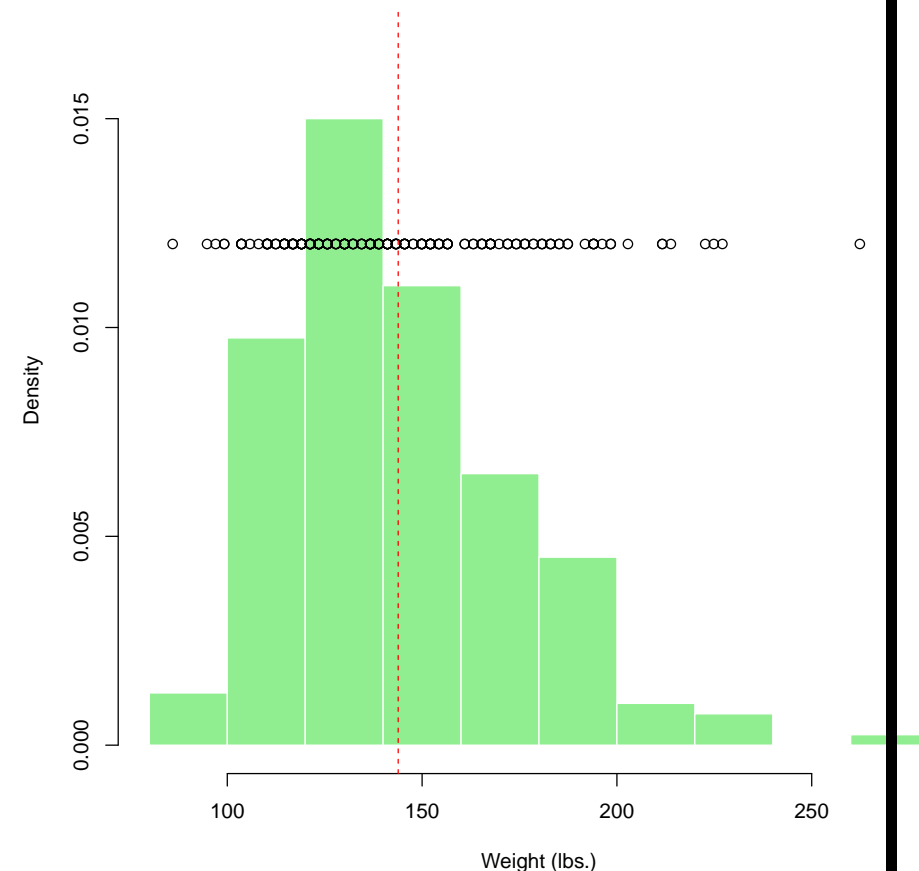
Consider a sample $\{y_1, \dots, y_n\} \longrightarrow$

Model : $Y_i = \beta_0 + \varepsilon_i$

with $\mathbb{E}(\varepsilon) = 0$ and $\text{Var}(\varepsilon) = \sigma^2$

ε is some unpredictable noise

$\hat{Y} = \bar{y}$ is our '*best guess*'...



Predicting means estimating $\mathbb{E}(Y)$.

Recall that $\mathbb{E}(Y) = \underset{y \in \mathbb{R}}{\operatorname{argmin}} \{ \|Y - y\|_{L_2} \} = \underbrace{\underset{y \in \mathbb{R}}{\operatorname{argmin}} \{ \mathbb{E}([Y - y]^2) \}}_{\text{least squares}}$

ε

Best estimate with some confidence

E.g. predicting someone's weight (Y)

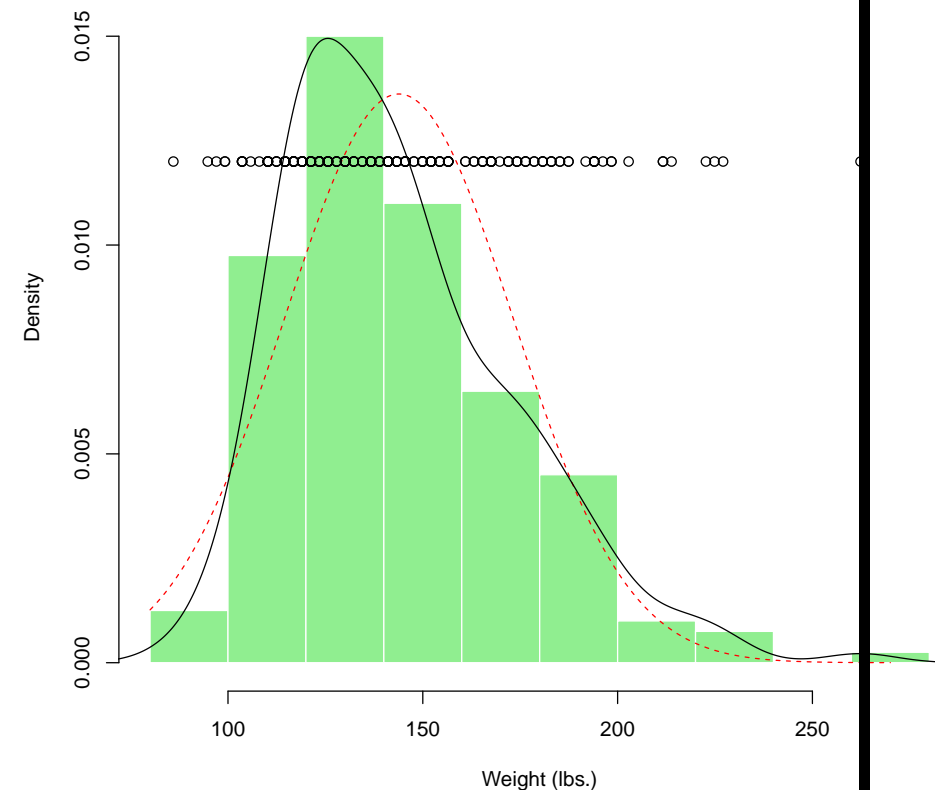
Give an **interval** $[y_-, y_+]$ such that

$$\mathbb{P}(Y \in [y_-, y_+]) = 1 - \alpha$$

Confidence intervals can be derived if

we can **estimate the distribution of Y**

$$F(y) = \mathbb{P}(Y \leq y) \text{ or } f(y) = \left. \frac{dF(x)}{dx} \right|_{x=y}$$



(related to the idea of “*quantifying uncertainty*” in our prediction...)

Parametric inference

E.g. predicting someone's weight (Y)

Assume that $F \in \mathcal{F} = \{F_{\theta}, \theta \in \Theta\}$

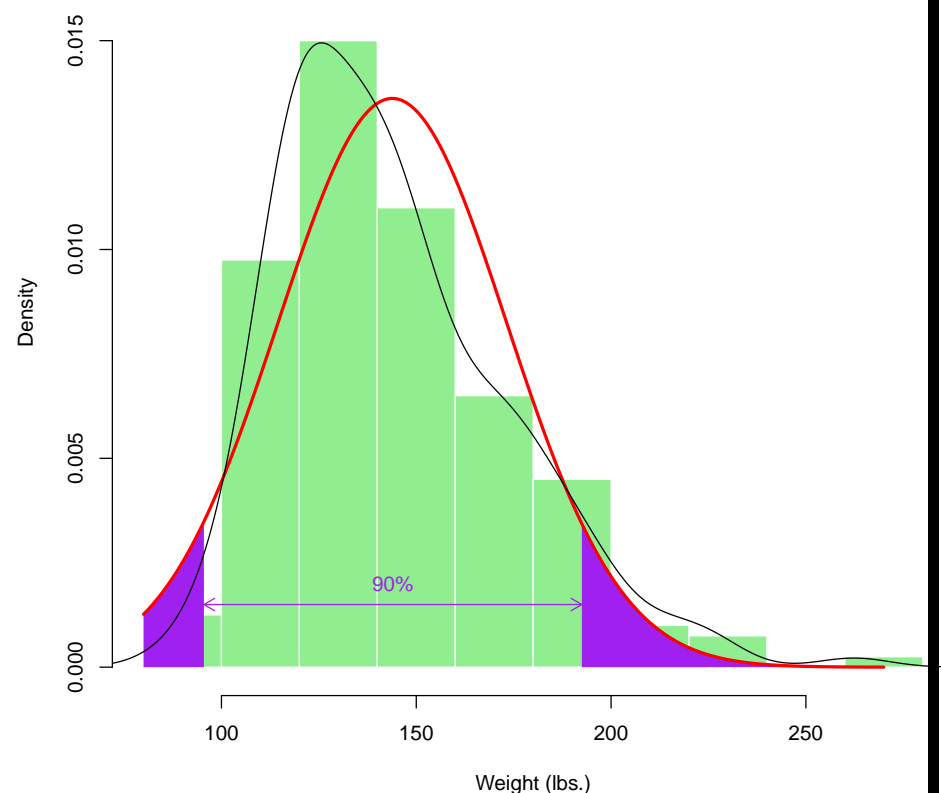
1. Provide an estimate $\hat{\theta}$
2. Compute bound estimates

$$\hat{y}_- = F_{\hat{\theta}}^{-1}(\alpha/2)$$

$$\hat{y}_+ = F_{\hat{\theta}}^{-1}(1 - \alpha/2)$$

Standard estimation technique :

→ maximum likelihood techniques



$$\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmax}} \left\{ \underbrace{\sum_{i=1}^n \log f_{\theta}(y_i)}_{\text{log likelihood}} \right\} \begin{cases} \text{explicit (analytical) expression for } \hat{\theta} \\ \text{numerical optimization (Newton Raphson)} \end{cases}$$

Non-parametric inference

E.g. predicting someone's weight (Y)

1. Empirical distribution function

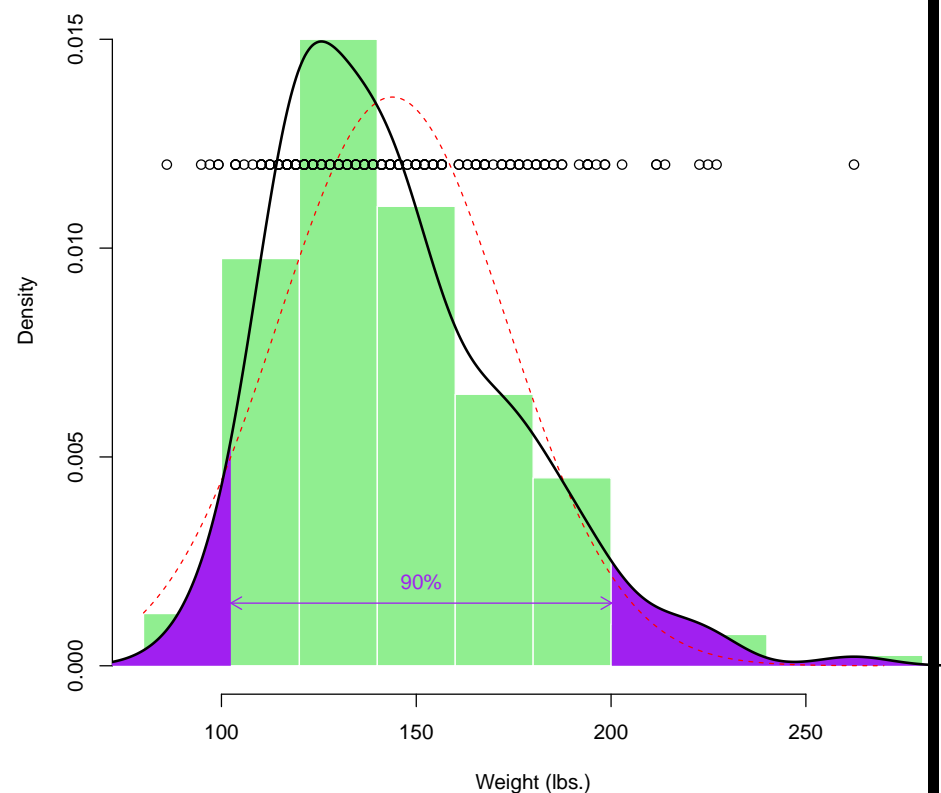
$$\hat{F}(y) = \frac{1}{n} \sum_{i=1}^n \underbrace{\mathbf{1}(y_i \leq y)}_{\#\{i \text{ such that } y_i \leq y\}}$$

natural estimator for $\mathbb{P}(Y \leq y)$

2. Compute bound estimates

$$\hat{y}_- = \hat{F}^{-1}(\alpha/2)$$

$$\hat{y}_+ = \hat{F}^{-1}(1 - \alpha/2)$$



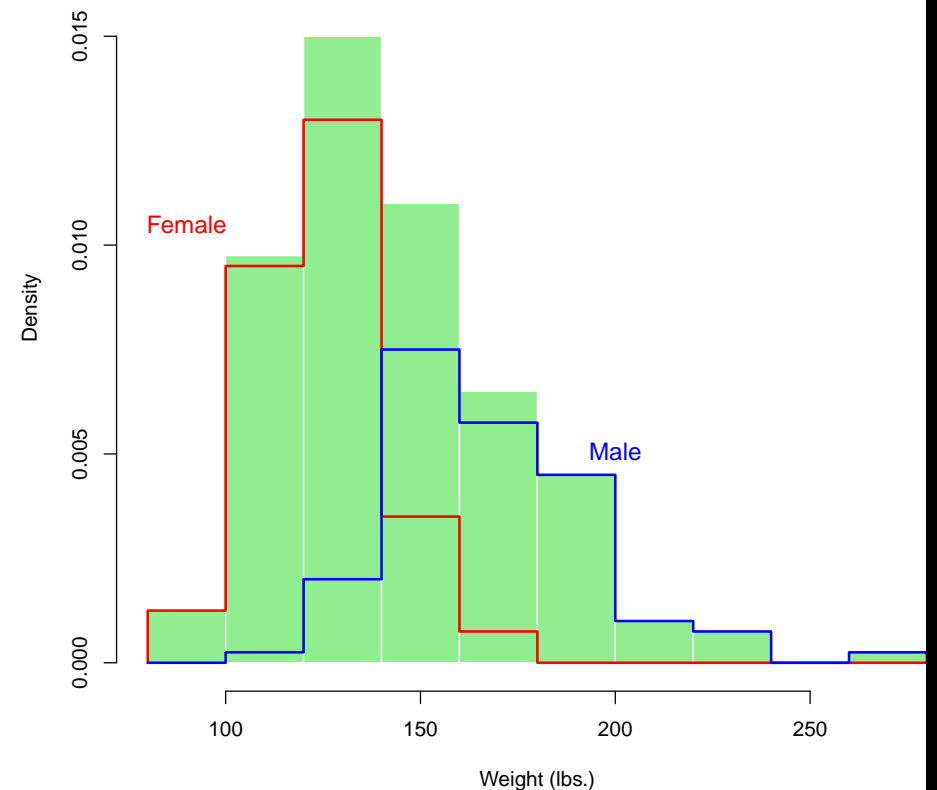
Prediction using some covariates

E.g. predicting someone's weight (Y)
based on his/her sex (X_1)

$$\text{Model : } Y_i = \begin{cases} \beta_F + \varepsilon_i & \text{if } X_{1,i} = F \\ \beta_H + \varepsilon_i & \text{if } X_{1,i} = M \end{cases}$$

$$\text{or } Y_i = \underbrace{\beta_0}_{\beta_M} + \underbrace{\beta_1}_{\beta_F - \beta_M} \mathbf{1}(X_{1,i} = F) + \varepsilon_i$$

with $\mathbb{E}(\varepsilon) = 0$ and $\text{Var}(\varepsilon) = \sigma^2$



Prediction using some (categorical) covariates

E.g. predicting someone's weight (Y)

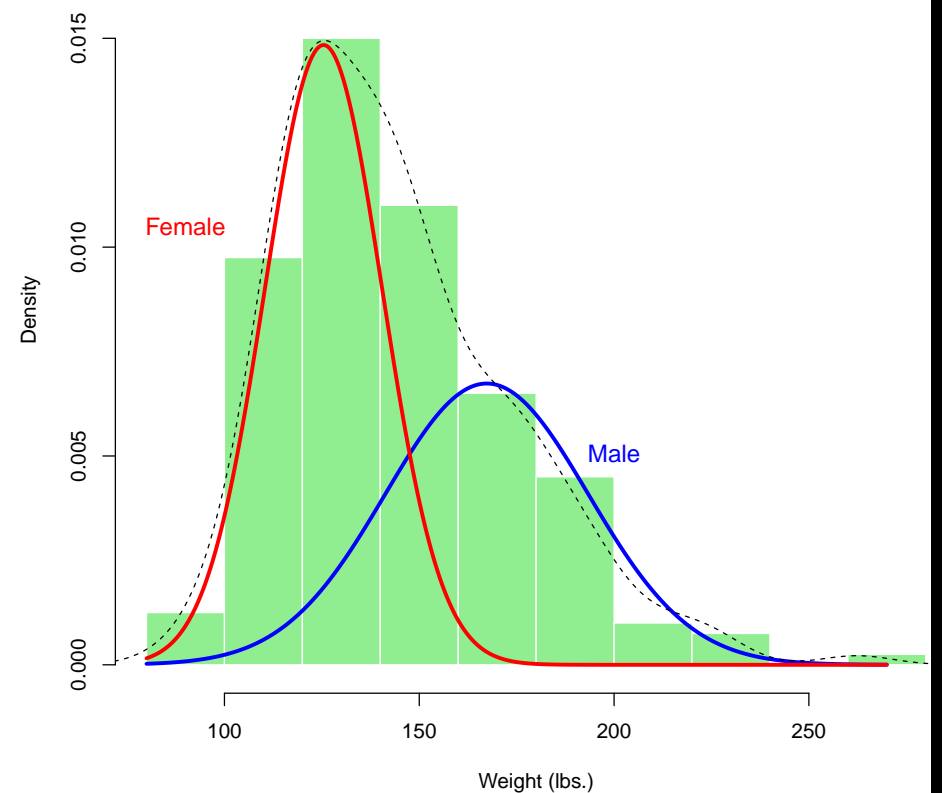
based on his/her sex (X_1)

Conditional parametric model

assume that $Y|X_1 = x_1 \sim F_{\theta(x_1)}$

$$\text{i.e. } Y_i \sim \begin{cases} F_{\theta_F} & \text{if } X_{1,i} = \text{F} \\ F_{\theta_M} & \text{if } X_{1,i} = \text{M} \end{cases}$$

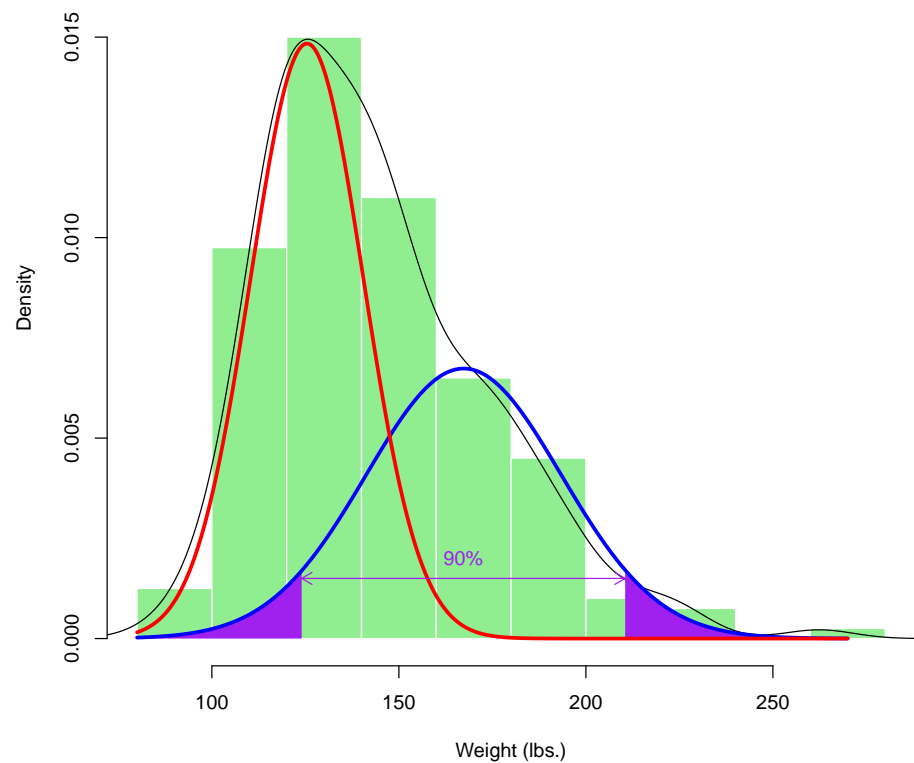
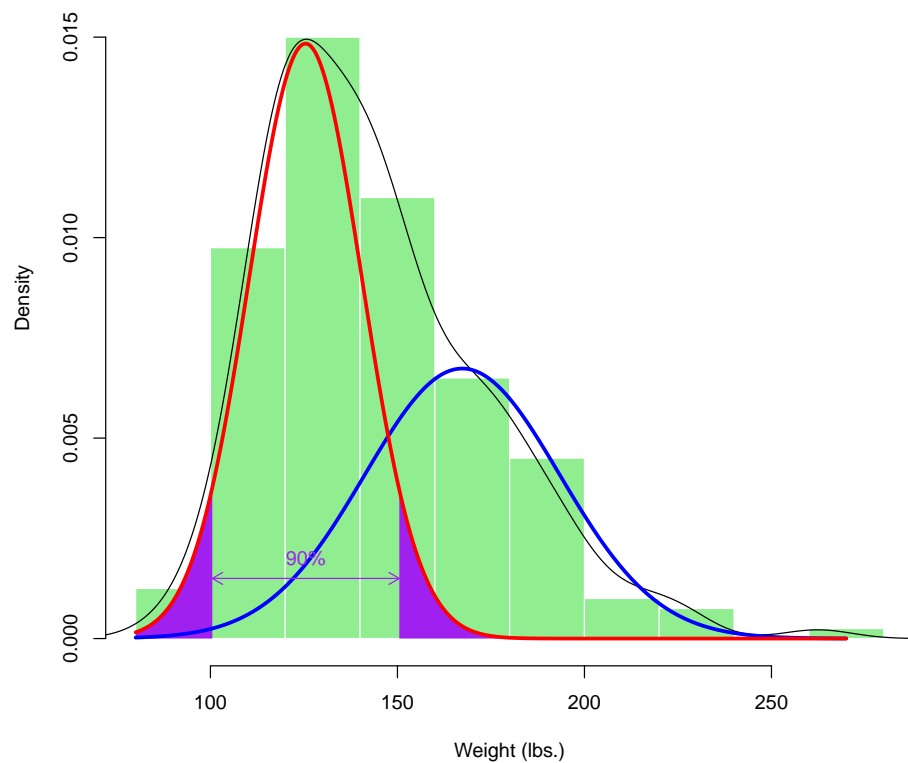
→ our prediction will be
conditional on the covariate



Prediction using some (categorical) covariates

Prediction of Y when $X_1 = F$

Prediction of Y when $X_1 = M$



Linear Models, and Ordinary Least Squares

E.g. predicting someone's weight (Y)

based on his/her height (X_2)

Linear Model : $Y_i = \beta_0 + \beta_2 X_{2,i} + \varepsilon_i$

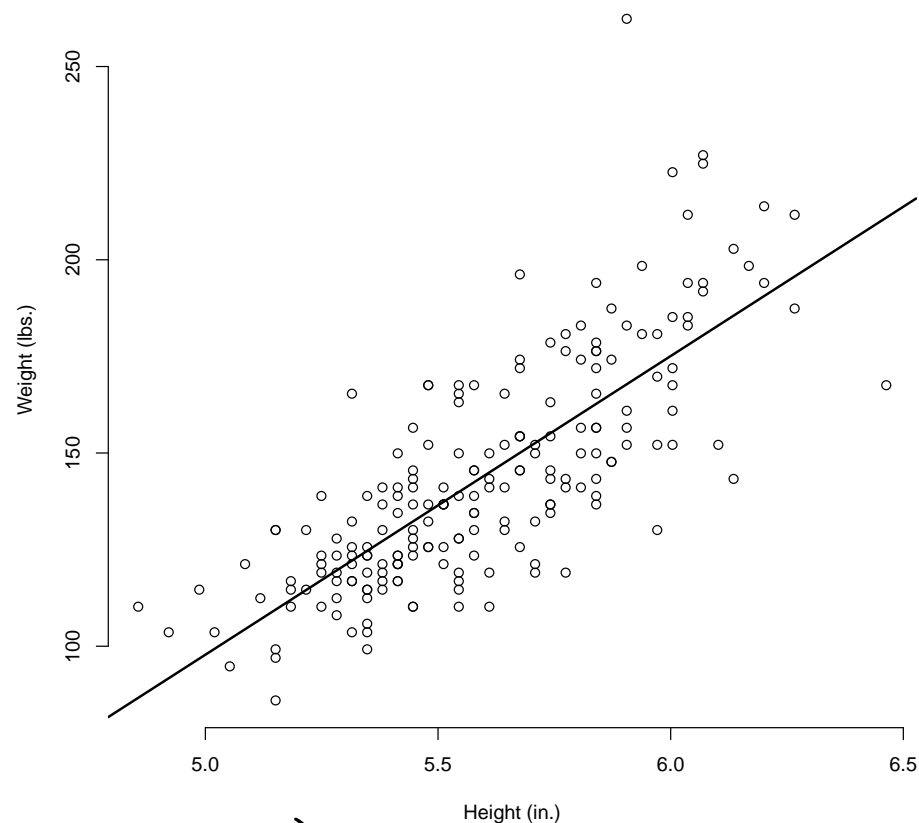
with $\mathbb{E}(\varepsilon) = 0$ and $\text{Var}(\varepsilon) = \sigma^2$

Conditional parametric model

assume that $Y|X_2 = x_2 \sim F_{\theta(x_2)}$

E.g. Gaussian Linear Model

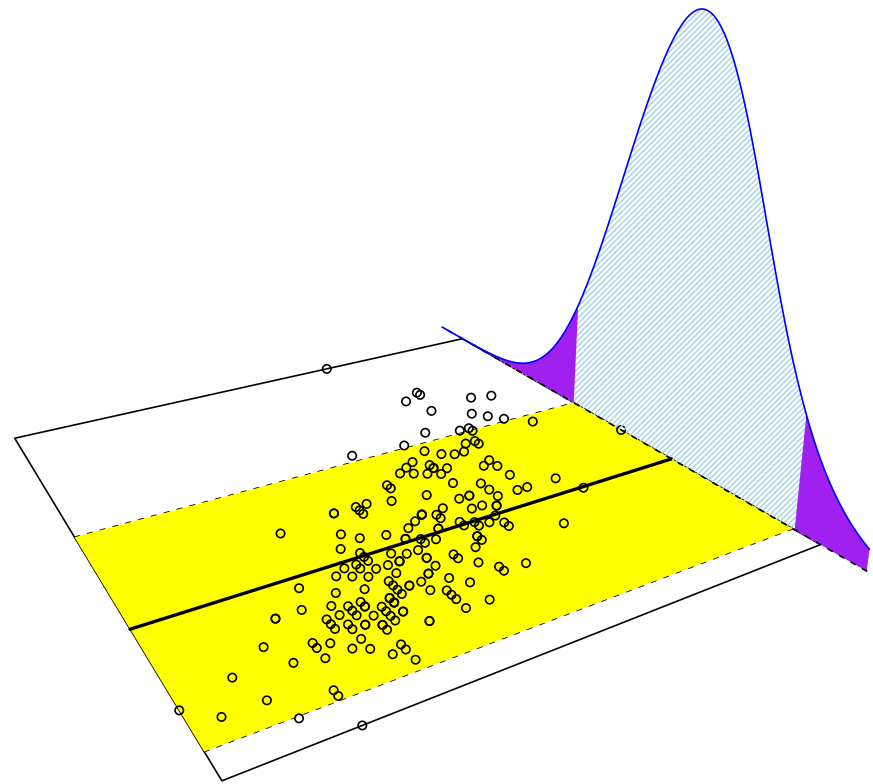
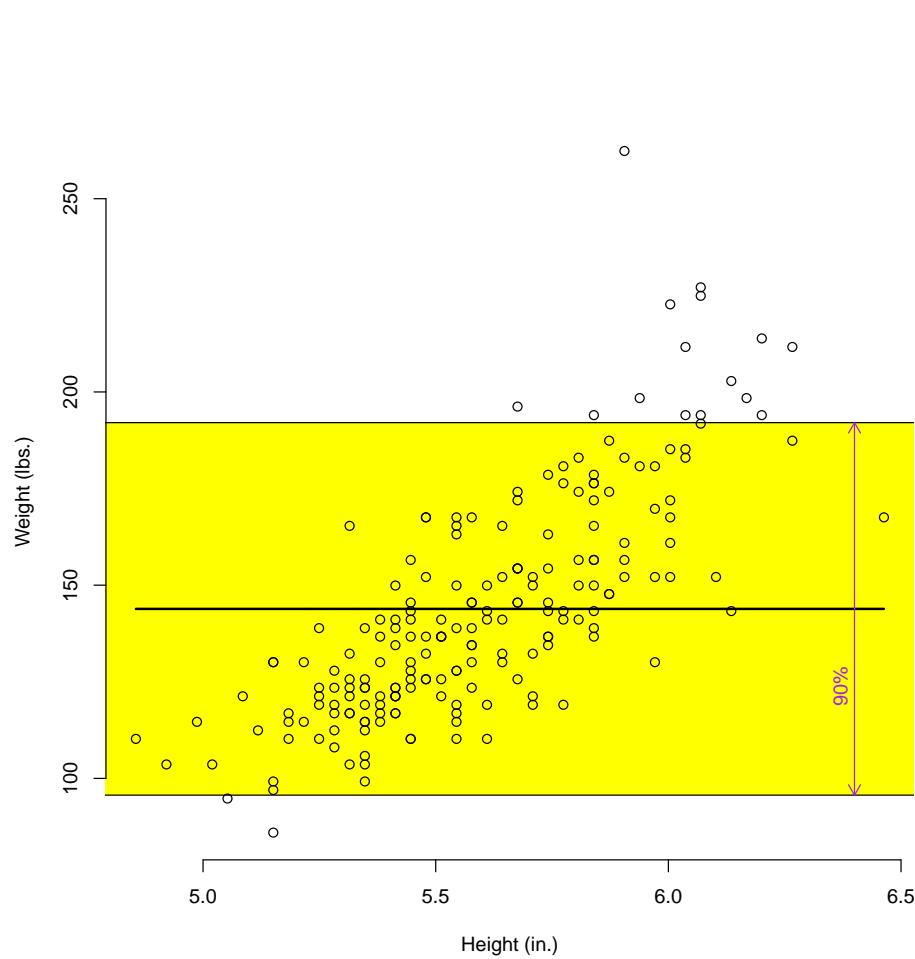
$Y|X_2 = x_2 \sim \mathcal{N}(\underbrace{\mu(x_2)}_{\beta_0 + \beta_2 x_2}, \underbrace{\sigma^2(x_2)}_{\sigma^2})$



→ ordinary least squares, $\hat{\beta} = \operatorname{argmin} \left\{ \sum_{i=1}^n [Y_i - \mathbf{X}_i^{\top} \beta]^2 \right\}$

$\hat{\beta}$ is also the M.L. estimator of β

Prediction using no covariates



Prediction using a categorical covariates

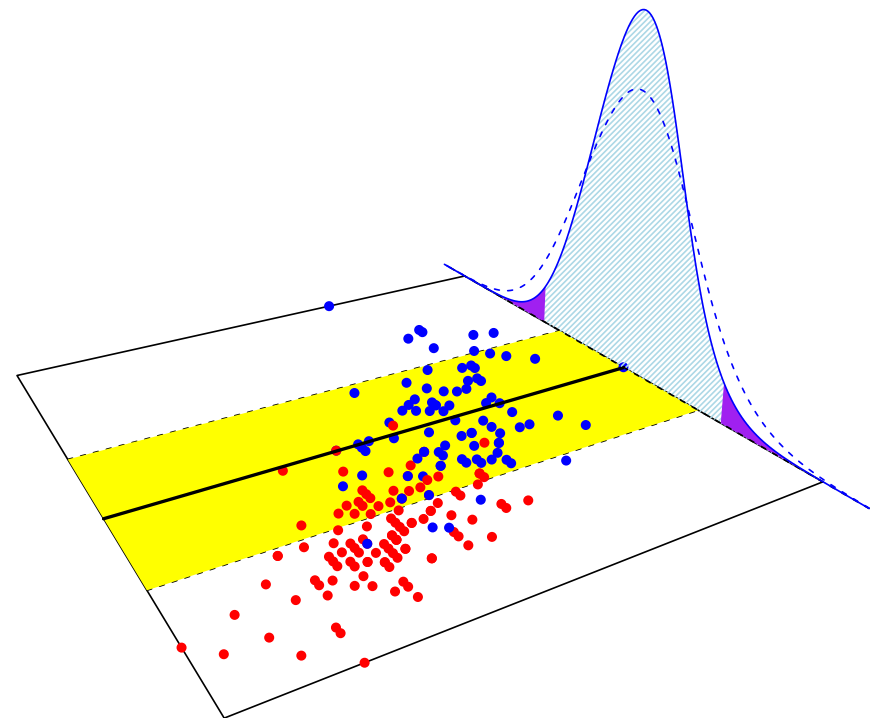
E.g. predicting someone's weight (Y)
based on his/her sex (X_1)

E.g. Gaussian linear model

$$Y|X_1 = M \sim \mathcal{N}(\mu_M, \sigma^2)$$

$$\hat{\mathbb{E}}(Y|X_1 = M) = \frac{1}{n_M} \sum_{i: X_{1,i}=M} Y_i = \hat{Y}(M)$$

$$Y \in \left[\hat{Y}(M) \pm \underbrace{u_{1-\alpha/2}}_{1.96} \cdot \hat{\sigma} \right]$$



Remark In the linear model, $\text{Var}(\varepsilon) = \sigma^2$ does not depend on X_1 .

Prediction using a categorical covariates

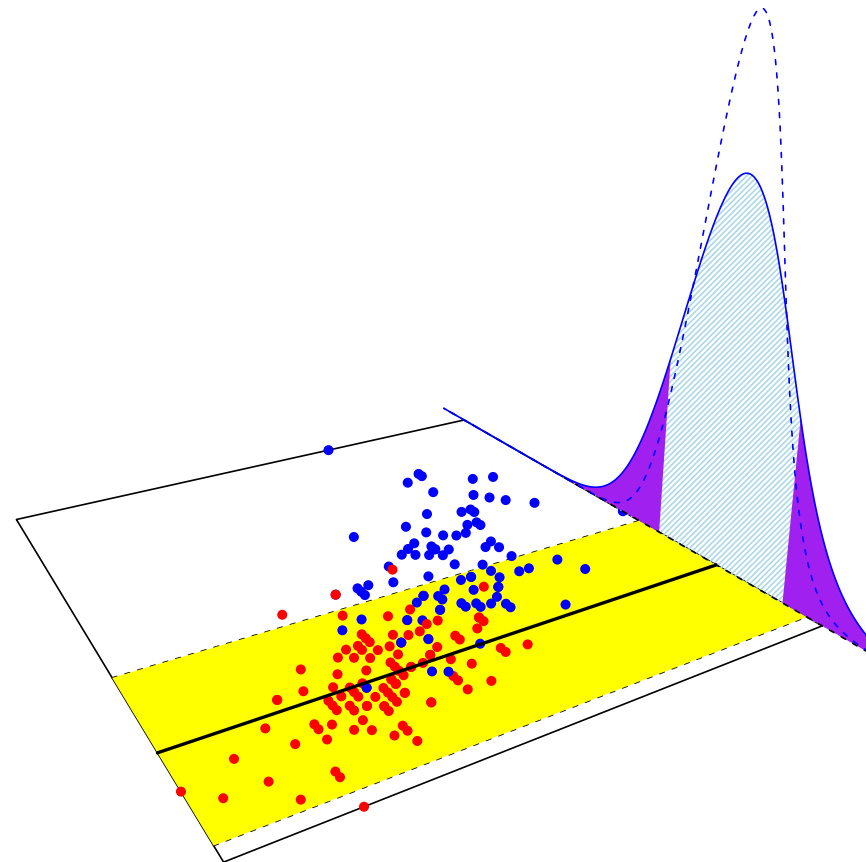
E.g. predicting someone's weight (Y)
based on his/her sex (X_1)

E.g. Gaussian linear model

$$Y | X_1 = F \sim \mathcal{N}(\mu_F, \sigma^2)$$

$$\hat{\mathbb{E}}(Y | X_1 = F) = \frac{1}{n_F} \sum_{i: X_{1,i} = F} Y_i = \hat{Y}(F)$$

$$Y \in [\hat{Y}(F) \pm \underbrace{u_{1-\alpha/2}}_{1.96} \cdot \hat{\sigma}]$$



Remark In the linear model, $\text{Var}(\varepsilon) = \sigma^2$ does not depend on X_1 .

Prediction using a continuous covariates

E.g. predicting someone's weight (Y)
based on his/her height (X_2)

E.g. Gaussian linear model

$$Y | X_2 = x_2 \sim \mathcal{N}(\beta_0 + \beta_1 x_2, \sigma^2)$$

$$\hat{\mathbb{E}}(Y | X_2 = x_2) = \hat{\beta}_0 + \hat{\beta}_1 x_2 = \hat{Y}(x_2)$$

$$Y \in [\hat{Y}(x_2) \pm \underbrace{u_{1-\alpha/2}}_{1.96} \cdot \hat{\sigma}]$$

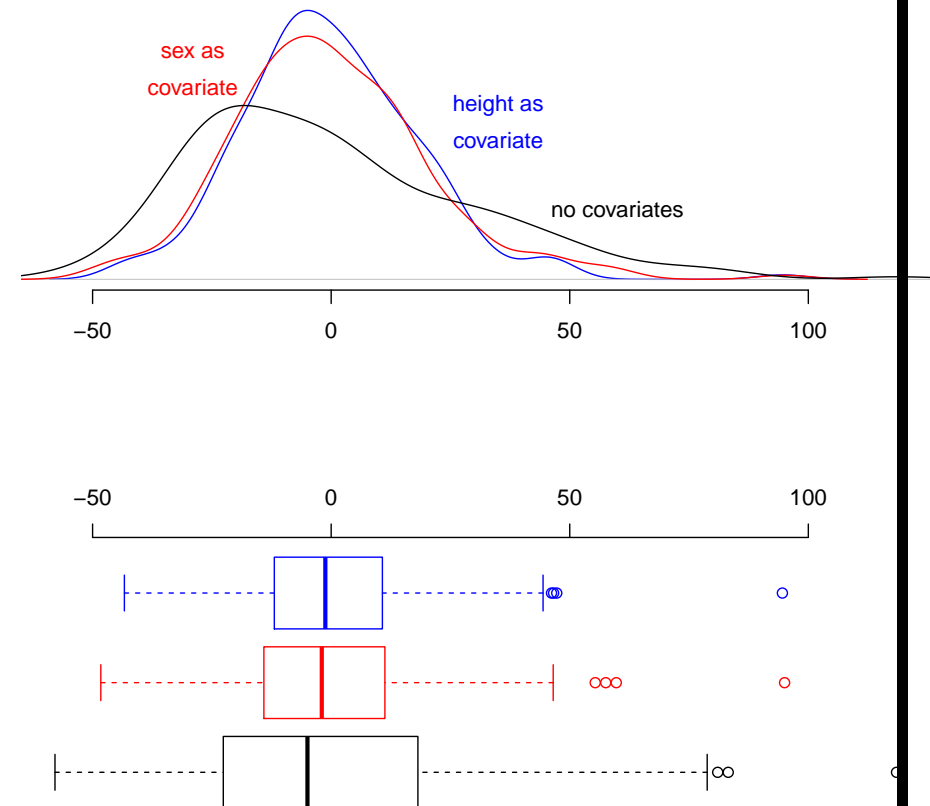
Improving our prediction ?

(Empirical) residuals, $\hat{\varepsilon}_i = Y_i - \underbrace{X_i^\top \hat{\beta}}_{\hat{Y}_i}$

R^2 or log-likelihood

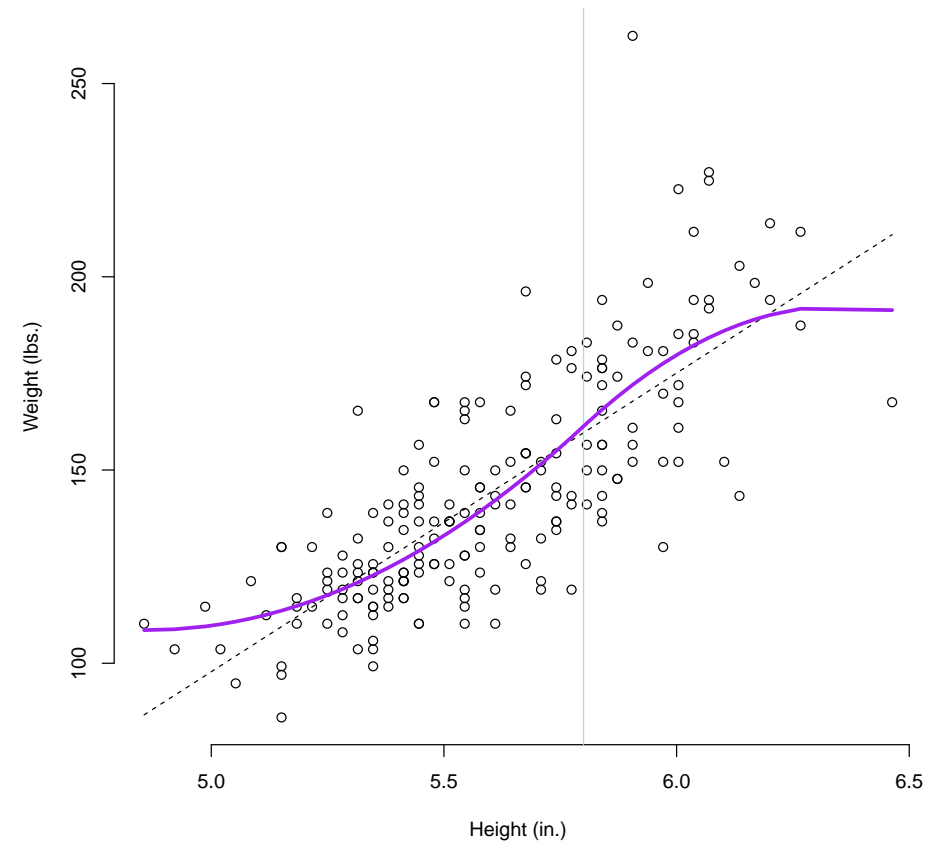
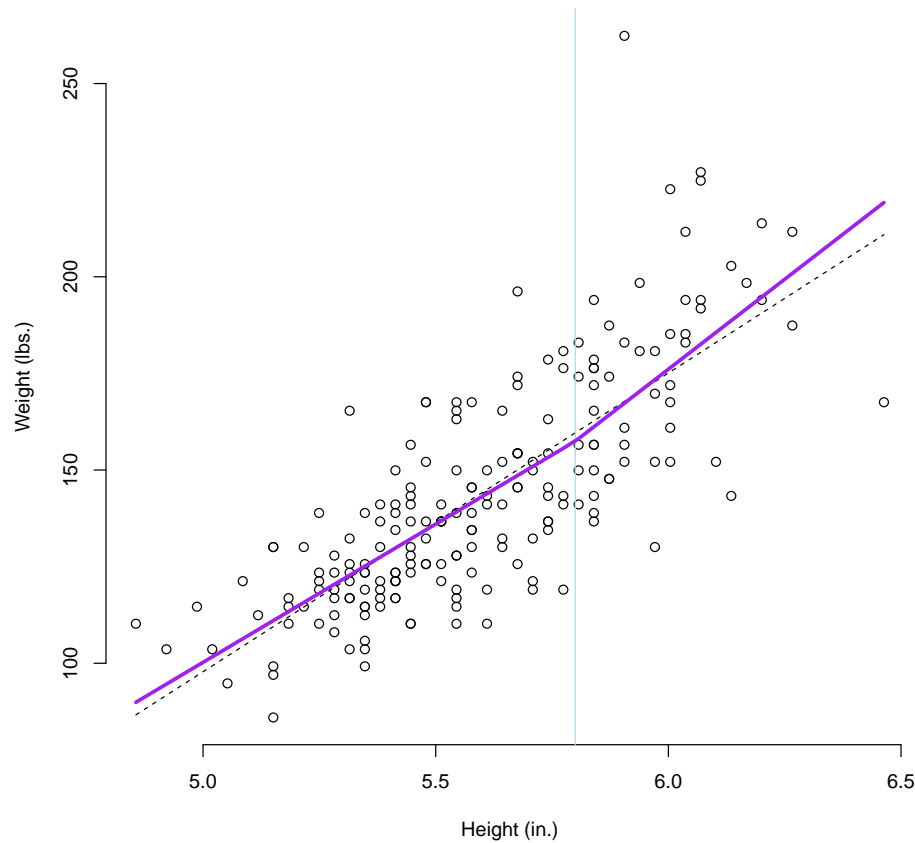
parsimony principle ?

→ penalizing the likelihood with
the number of covariates
Akaike (AIC) or
Schwarz (BIC) criteria



Relaxing the linear assumption in predictions

Use of *b*-spline function basis to estimate $\mu(\cdot)$ where $\mu(\mathbf{x}) = \mathbb{E}(Y|\mathbf{X} = \mathbf{x})$



Relaxing the linear assumption in predictions

E.g. predicting someone's weight (Y)
based on his/her height (X_2)

E.g. Gaussian linear model

$$Y|X_2 = x_2 \sim \mathcal{N}(\mu(x_2), \sigma^2)$$

$$\hat{\mathbb{E}}(Y|X_2 = x_2) = \hat{\mu}(x_2) = \hat{Y}(x_2)$$

$$Y \in [\hat{Y}(x_2) \pm \underbrace{u_{1-\alpha/2}}_{1.96} \cdot \hat{\sigma}]$$

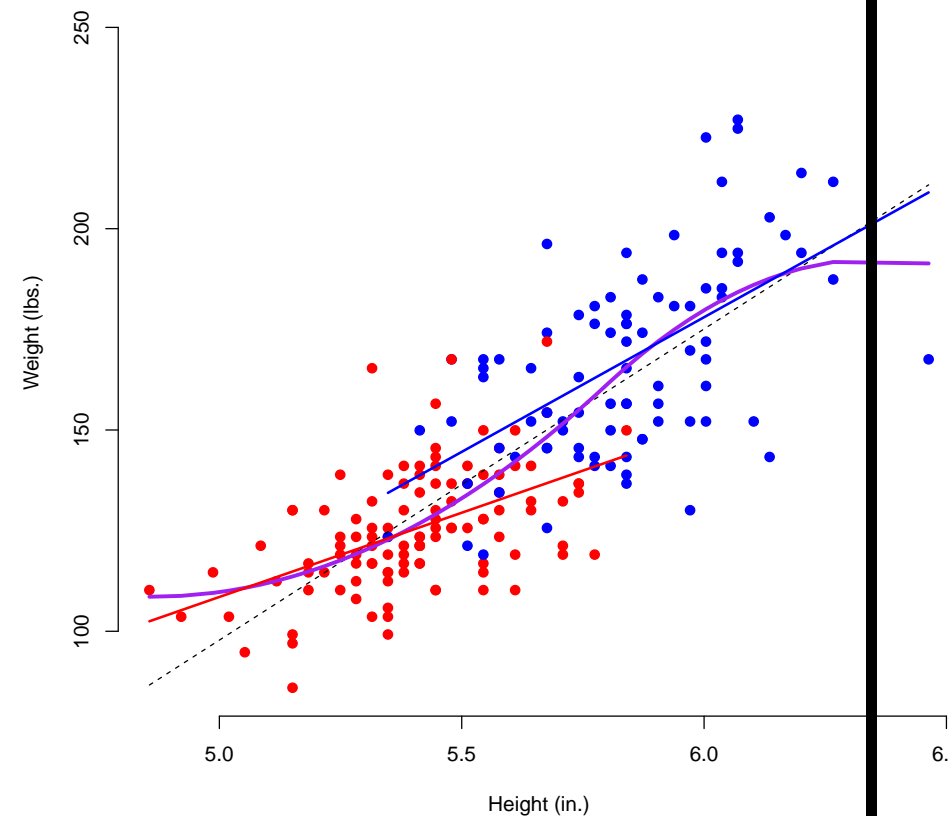
Gaussian model : $\mathbb{E}(Y|\mathbf{X} = \mathbf{x}) = \mu(\mathbf{x})$ (e.g. $\mathbf{x}^\top \boldsymbol{\beta}$) and $\text{Var}(Y|\mathbf{X} = \mathbf{x}) = \sigma^2$.

Nonlinearities and missing covariates

E.g. predicting someone's weight (Y)
 based on his/her height **and** sex
 → nonlinearities can be related to
 model misspecification

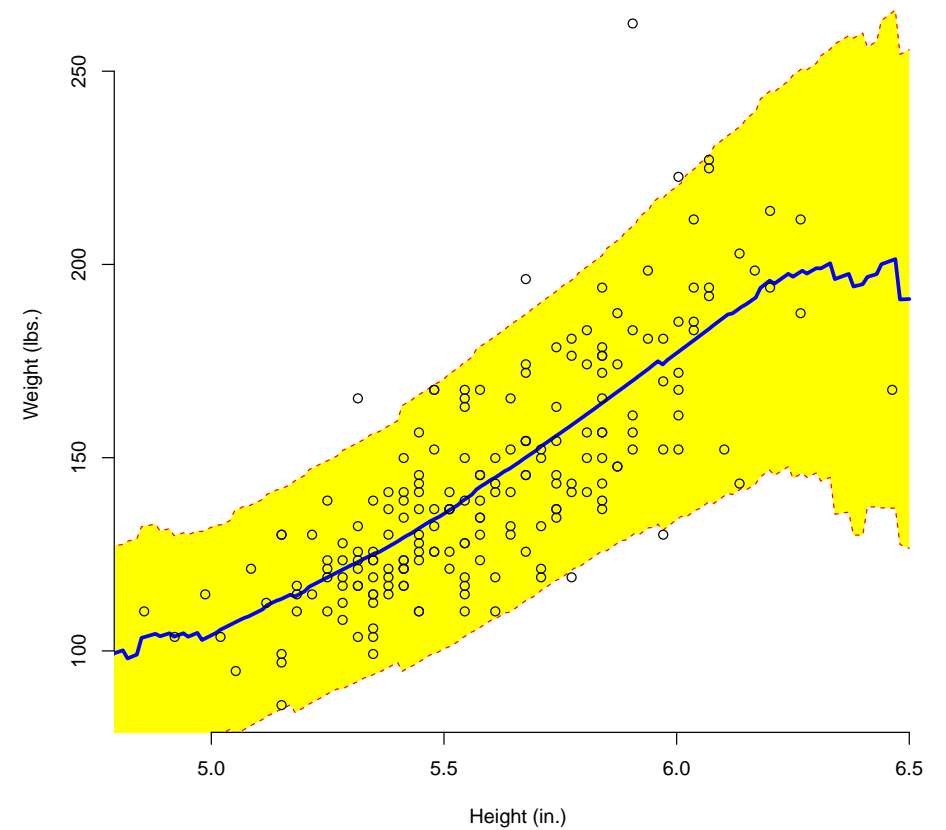
E.g. Gaussian linear model

$$Y_i = \begin{cases} \beta_{0,F} + \beta_{2,F} X_{2,i} + \varepsilon_i & \text{if } X_{1,i} = F \\ \beta_{0,M} + \beta_{2,M} X_{2,i} + \varepsilon_i & \text{if } X_{1,i} = M \end{cases}$$

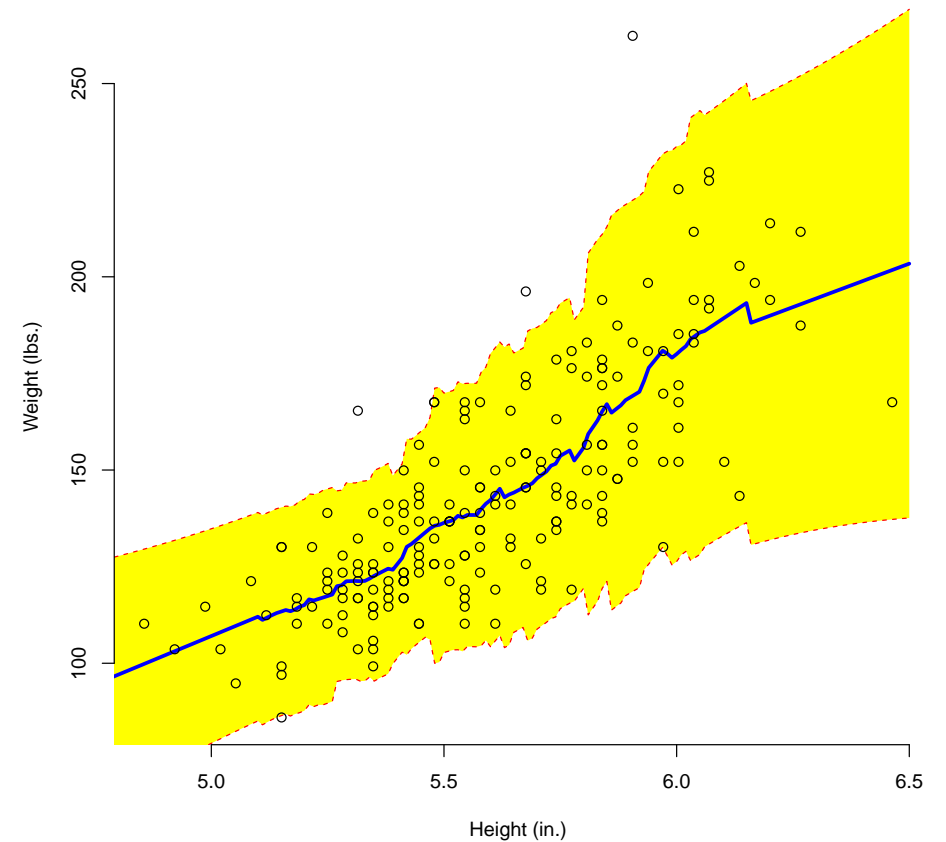


→ local linear regression, $\hat{\beta}_x = \operatorname{argmin} \left\{ \sum_{i=1}^n \omega_i(\mathbf{x}) \cdot [Y_i - \mathbf{X}_i^\top \beta]^2 \right\}$
 and set $\hat{Y}(\mathbf{x}) = \mathbf{x}^\top \hat{\beta}_x$

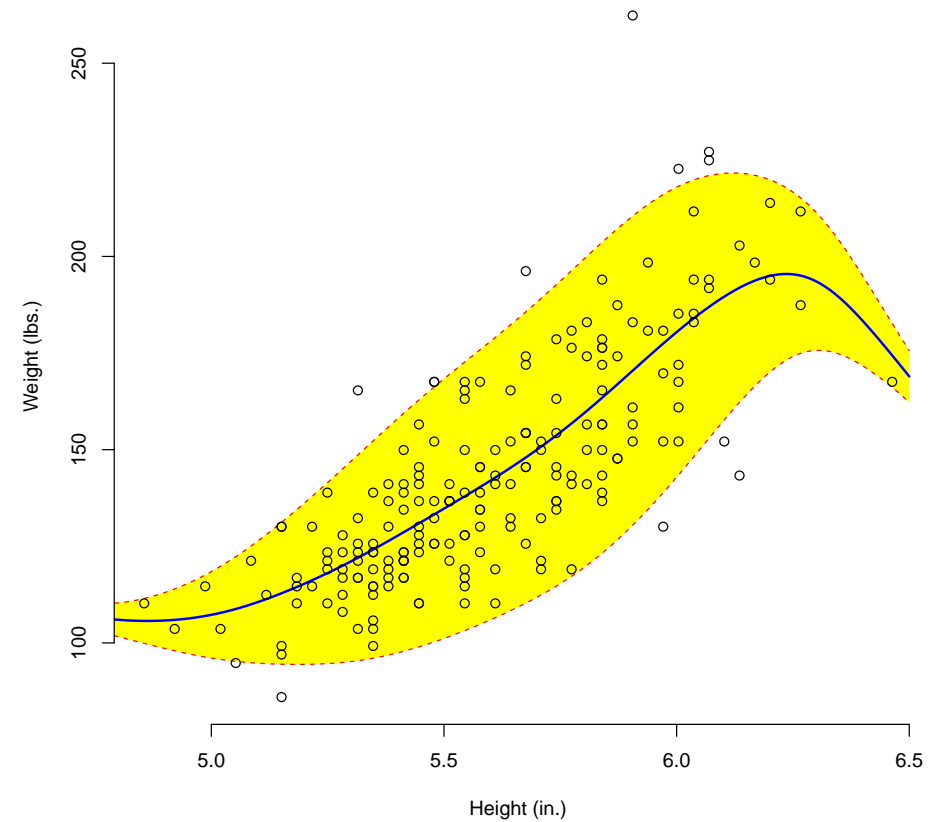
Local regression and smoothing techniques



k-nearest neighbours and smoothing techniques



Kernel regression and smoothing techniques



Multiple linear regression

E.g. predicting someone's misperception
of his/her weight (Y)

based on his/her $\underbrace{\text{height}}_{X_2}$ **and** $\underbrace{\text{weight}}_{X_3}$

→ **linear model**

$$\mathbb{E}(Y|X_2, X_3) = \beta_0 + \beta_2 X_2 + \beta_3 X_3$$

$$\text{Var}(Y|X_2, X_3) = \sigma^2$$

Multiple non-linear regression

E.g. predicting someone's misperception
of his/her weight (Y)

based on his/her $\underbrace{\text{height}}_{X_2}$ **and** $\underbrace{\text{weight}}_{X_3}$

→ **non-linear model**

$$\mathbb{E}(Y|X_2, X_3) = h(X_2, X_3)$$

$$\text{Var}(Y|X_2, X_3) = \sigma^2$$

Away from the Gaussian model

Y is not necessarily Gaussian

Y can be a counting variable,

E.g. Poisson $Y \sim \mathcal{P}(\lambda(\mathbf{x}))$

Y can be a FALSE-TRUE variable,

E.g. Binomial $Y \sim \mathcal{B}(p(\mathbf{x}))$

(see next section)

→ Generalized Linear Model

E.g. $Y|X_2 = x_2 \sim \mathcal{P}(e^{\beta_0 + \beta_2 x_2})$

Remark With a Poisson model, $\mathbb{E}(Y|X_2 = x_2) = \text{Var}(Y|X_2 = x_2)$.

Logistic regression

E.g. predicting someone's misperception
of his/her weight (Y)

$$Y_i = \begin{cases} 1 & \text{if prediction} > \text{observed weight} \\ 0 & \text{if prediction} \leq \text{observed weight} \end{cases}$$

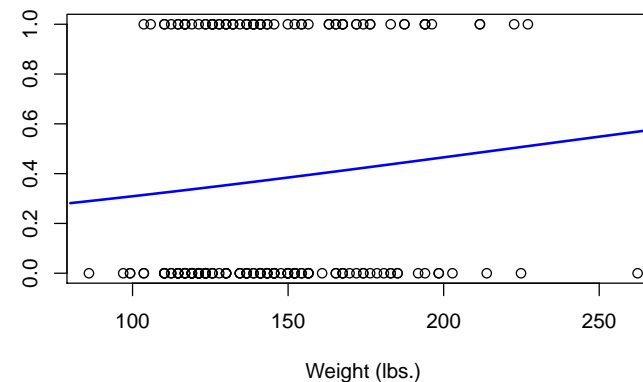
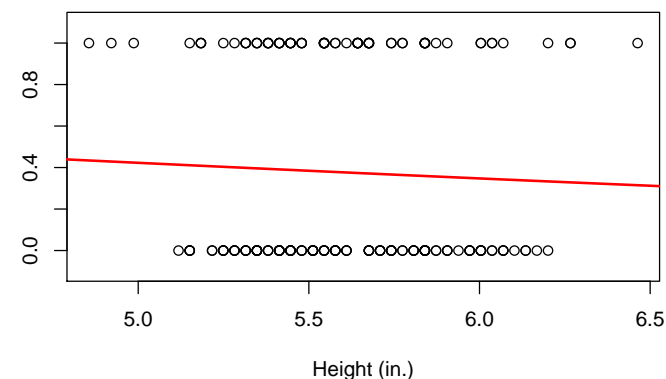
Bernoulli variable,

$$\mathbb{P}(Y = y) = p^y(1 - p)^{1-y}, \text{ where } y \in \{0, 1\}$$

→ logistic regression

$$\mathbb{P}(Y = y | \mathbf{X} = \mathbf{x}) = p(\mathbf{x})^y(1 - p(\mathbf{x}))^{1-y},$$

where $y \in \{0, 1\}$



Logistic regression

→ logistic regression

$$\mathbb{P}(Y = y | \mathbf{X} = \mathbf{x}) = p(\mathbf{x})^y (1 - p(\mathbf{x}))^{1-y},$$

where $y \in \{0, 1\}$

$$\text{Odds ratio } \frac{\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})}{\mathbb{P}(Y = 0 | \mathbf{X} = \mathbf{x})} = \exp(\mathbf{x}^\top \boldsymbol{\beta})$$

$$\mathbb{E}(Y | \mathbf{X} = \mathbf{x}) = \frac{\exp(\mathbf{x}^\top \boldsymbol{\beta})}{1 + \exp(\mathbf{x}^\top \boldsymbol{\beta})}$$

Estimation of $\boldsymbol{\beta}$?

→ maximum likelihood $\hat{\boldsymbol{\beta}}$ (Newton - Raphson)

Smoothed logistic regression

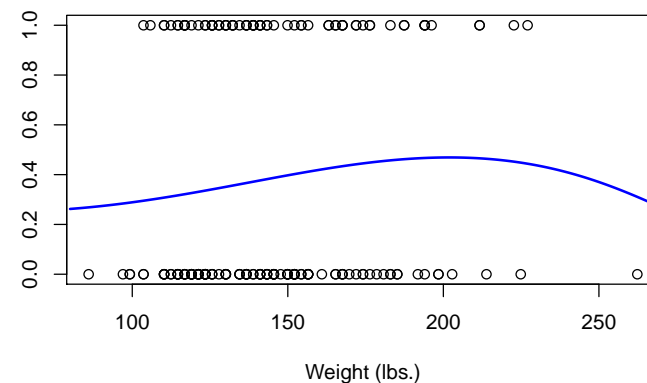
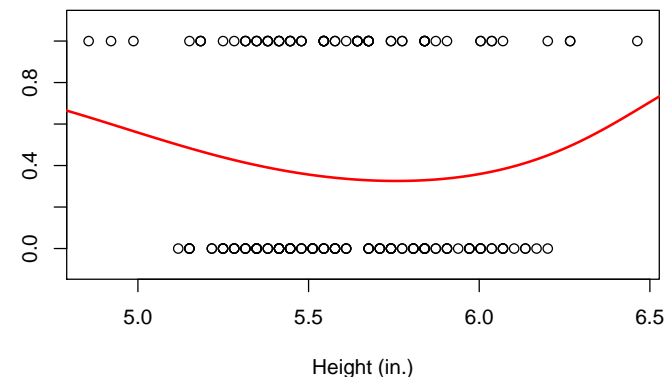
GLMs are linear since

$$\frac{\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})}{\mathbb{P}(Y = 0 | \mathbf{X} = \mathbf{x})} = \exp(\mathbf{x}^\top \boldsymbol{\beta})$$

—→ smooth nonlinear function instead

$$\frac{\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})}{\mathbb{P}(Y = 0 | \mathbf{X} = \mathbf{x})} = \exp(h(\mathbf{x}))$$

$$\mathbb{E}(Y | \mathbf{X} = \mathbf{x}) = \frac{\exp(h(\mathbf{x}))}{1 + \exp(h(\mathbf{x}))}$$



Smoothed logistic regression

—→ non linear logistic regression

$$\frac{\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})}{\mathbb{P}(Y = 0 | \mathbf{X} = \mathbf{x})} = \exp(h(\mathbf{x}))$$

$$\mathbb{E}(Y | \mathbf{X} = \mathbf{x}) = \frac{\exp(h(\mathbf{x}))}{1 + \exp(h(\mathbf{x}))}$$

Remark we do not predict Y here,
but $\mathbb{E}(Y | \mathbf{X} = \mathbf{x})$.

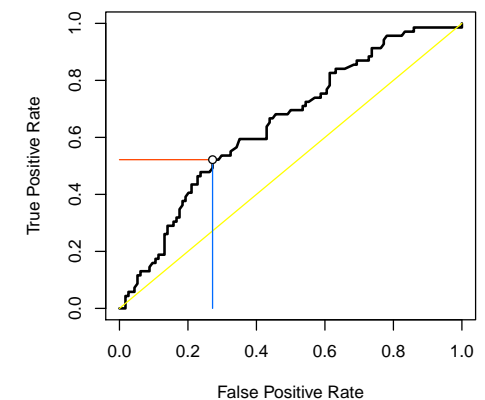
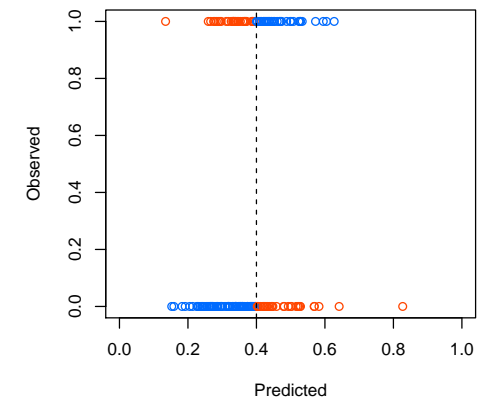
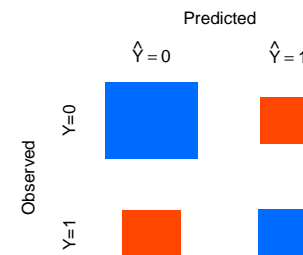
Predictive modeling for a $\{0, 1\}$ variable

What is a good $\{0, 1\}$ -model?

→ decision theory

$$\begin{cases} \text{if } \mathbb{P}(Y|\mathbf{X} = \mathbf{x}) \leq \mathbf{s}, \text{ then } \hat{Y} = 0 \\ \text{if } \mathbb{P}(Y|\mathbf{X} = \mathbf{x}) > \mathbf{s}, \text{ then } \hat{Y} = 1 \end{cases}$$

	$\hat{Y} = 0$	$\hat{Y} = 1$
$Y = 0$	fine	error
$Y = 1$	error	fine



R.O.C. curve

True positive rate

$$\begin{aligned} TP(s) &= \mathbb{P}(\hat{Y}(s) = 1 | Y = 1) \\ &= \frac{n_{\hat{Y}=1, Y=1}}{n_{Y=1}} \end{aligned}$$

False positive rate

$$\begin{aligned} FP(s) &= \mathbb{P}(\hat{Y}(s) = 1 | Y = 0) \\ &= \frac{n_{\hat{Y}=1, Y=0}}{n_{Y=0}} \end{aligned}$$

R.O.C. curve is

$$\{FP(s), TP(s), s \in (0, 1)\}$$

(see also model gain curve)

Classification tree (CART)

If Y is a TRUE-FALSE variable
prediction is a **classification** problem.

$$\mathbb{E}(Y|\mathbf{X} = \mathbf{x}) = p_j \text{ if } \mathbf{x} \in A_j$$

where A_1, \dots, A_k are disjoint
regions of the \mathbf{X} -space.

Classification tree (CART)

→ iterative process

Step 1. Find two subset of indices

either $A_1 = \{i, X_{1,i} < s\}$

and $A_2 = \{i, X_{1,i} > s\}$

or $A_1 = \{i, X_{2,i} < s\}$

and $A_2 = \{i, X_{2,i} > s\}$

maximize homogeneity within subsets

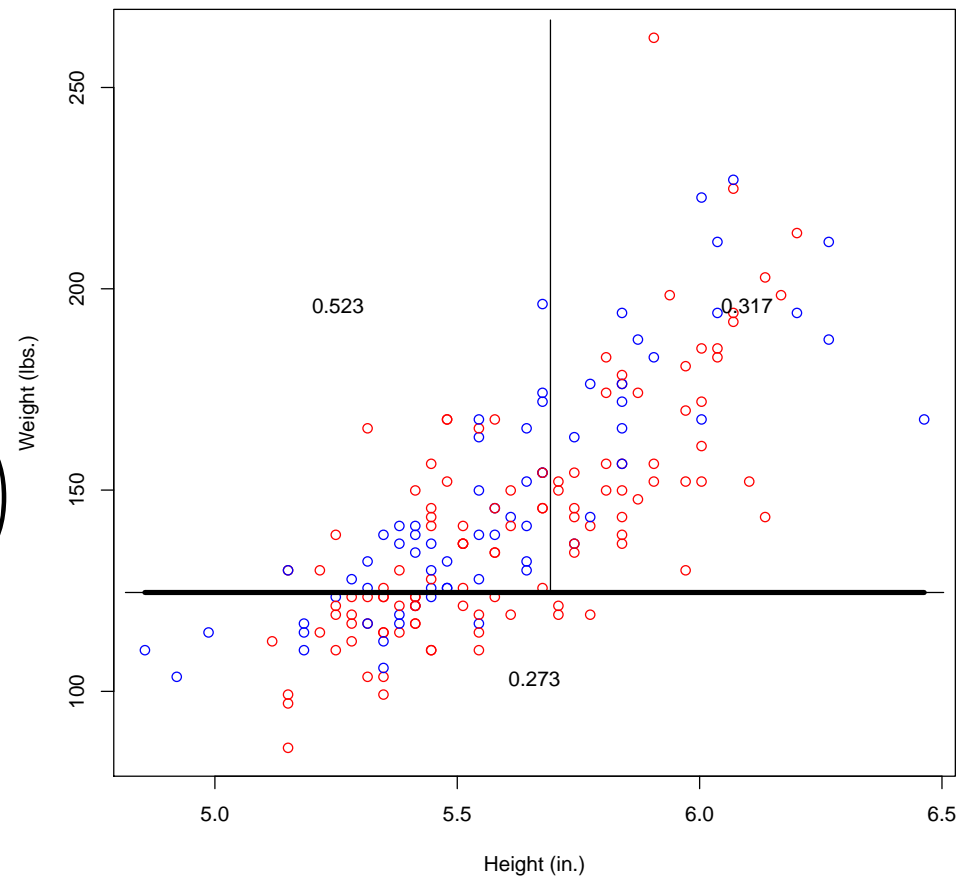
& maximize heterogeneity between subsets

Classification tree (CART)

Need an impurity criteria

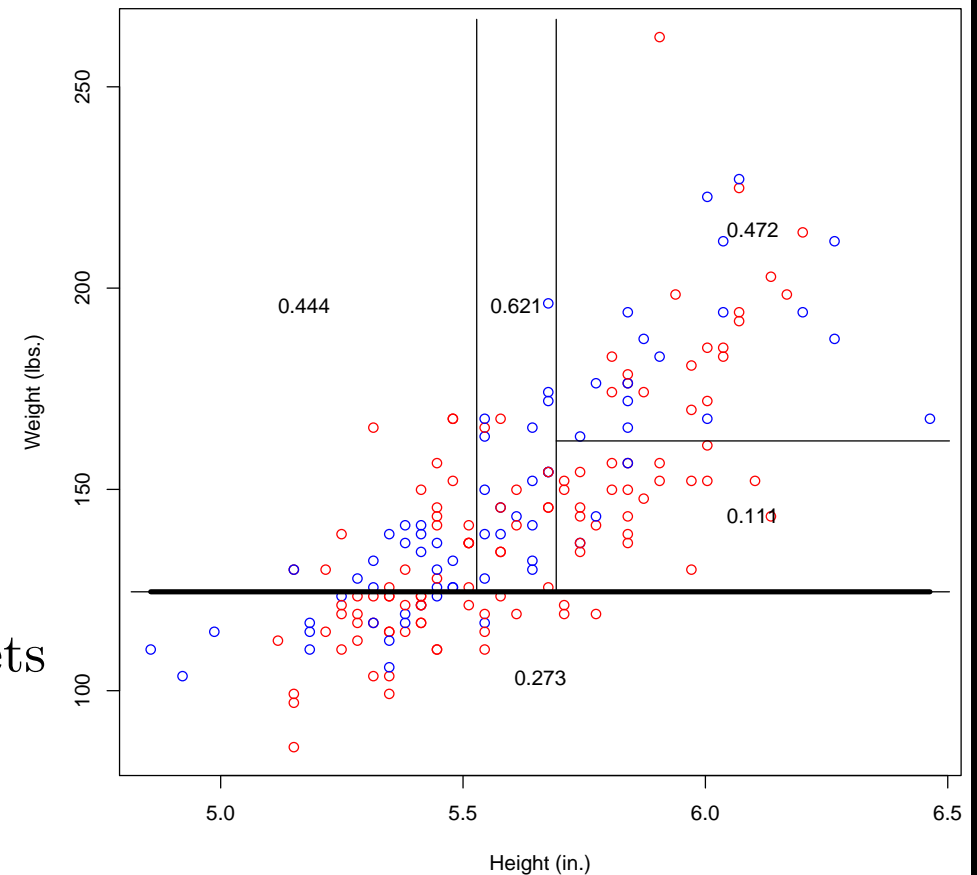
E.g. Gini index

$$- \sum_{\mathbf{x} \in \{A_1, A_2\}} \frac{n_{\mathbf{x}}}{n} \sum_{y \in \{0,1\}} \frac{n_{\mathbf{x},y}}{n_{\mathbf{x}}} \cdot \left(1 - \frac{n_{\mathbf{x},y}}{n_{\mathbf{x}}}\right)$$

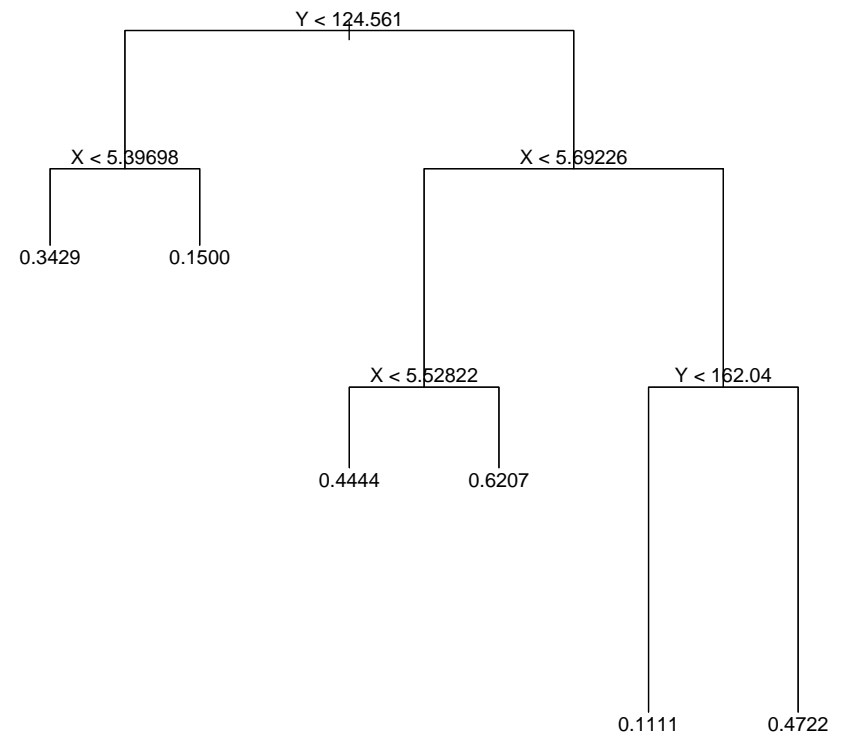
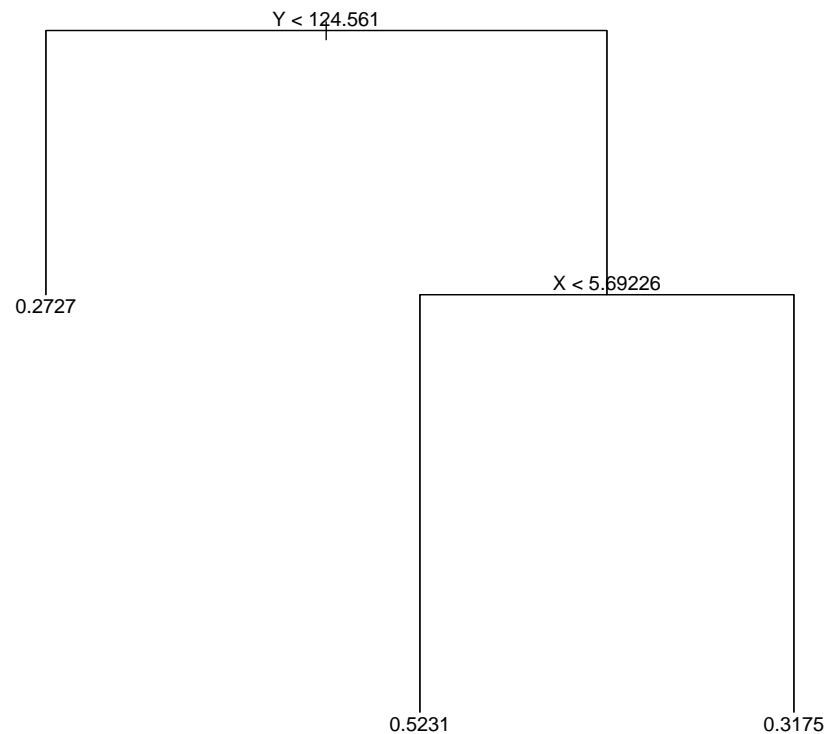


Classification tree (CART)

Step k . Given partition A_1, \dots, A_k
find which subset A_j will be divided,
either according to X_1
or according to X_2
maximize homogeneity within subsets
& maximize heterogeneity between subsets



Visualizing classification trees (CART)



From trees to forests

Problem CART tree are not robust

—→ **boosting** and **bagging**

use bootstrap : resample in the data

and generate a classification tree

repeat this resampling strategy

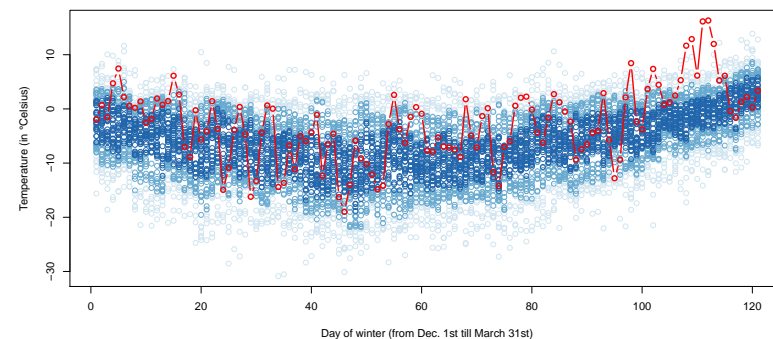
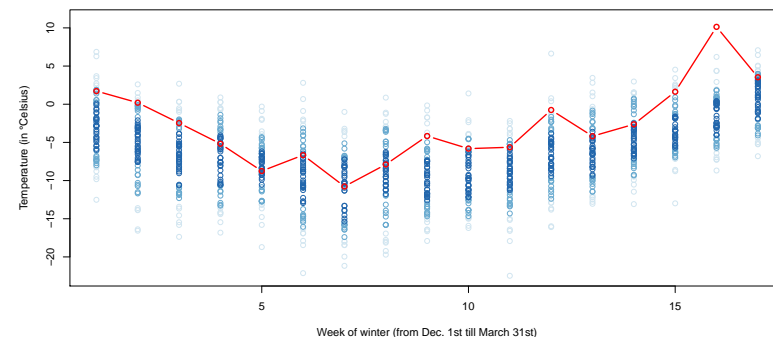
Then aggregate all the trees

A short word on functional data

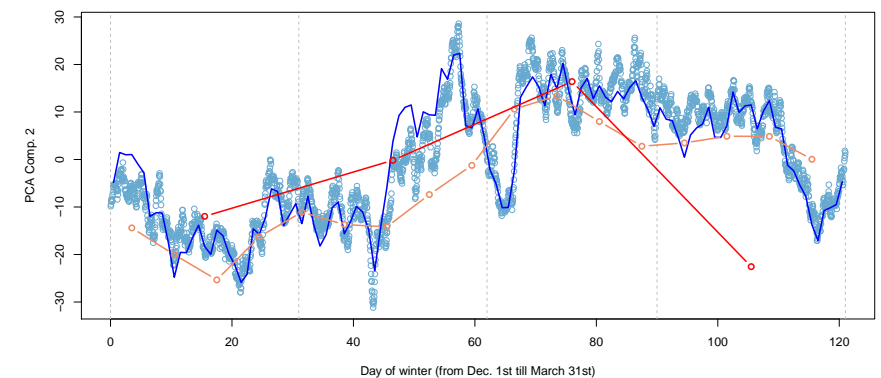
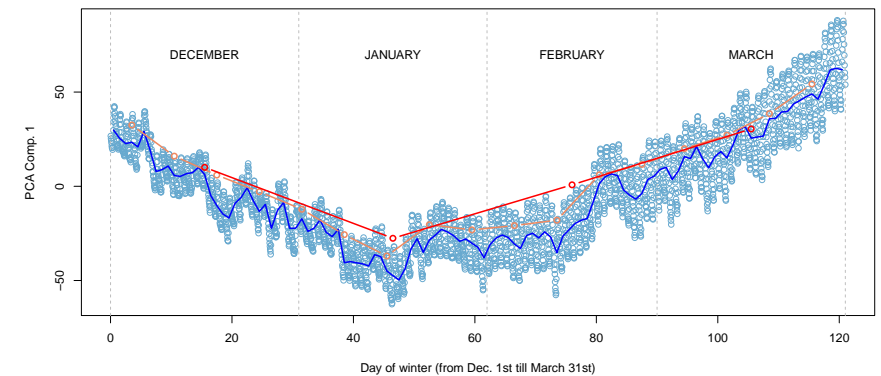
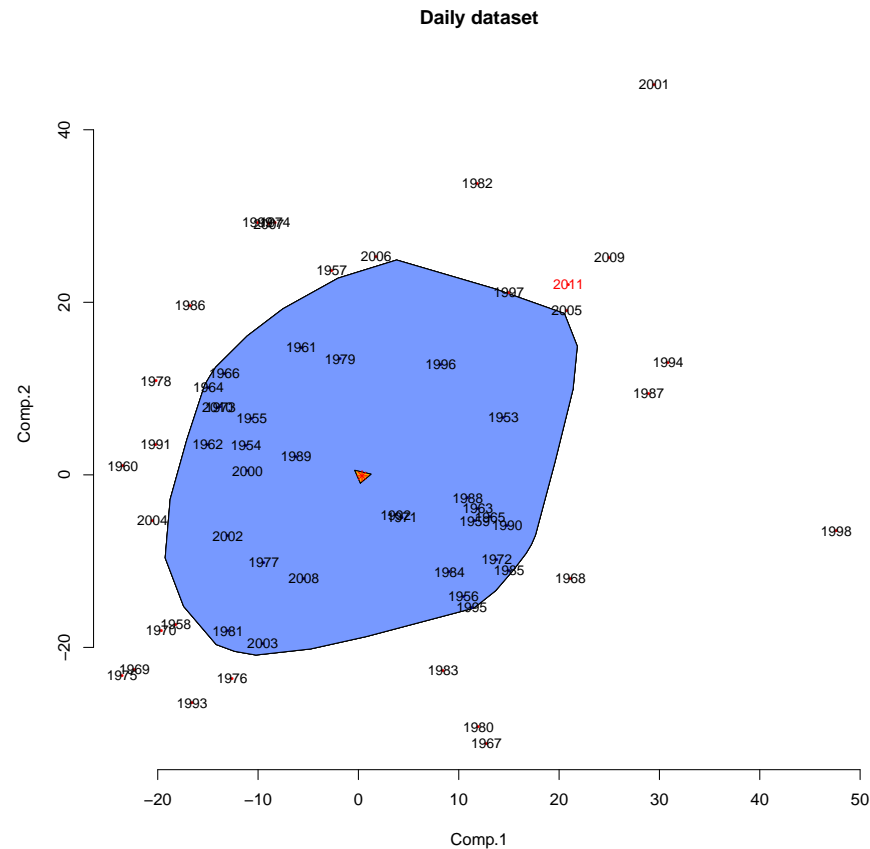
Individual data, $\{Y_i, (X_{1,i}, \dots, X_{k,i}, \dots)\}$

Functional data, $\{\mathbf{Y}_i = (Y_{i,1}, \dots, Y_{i,t}, \dots)\}$

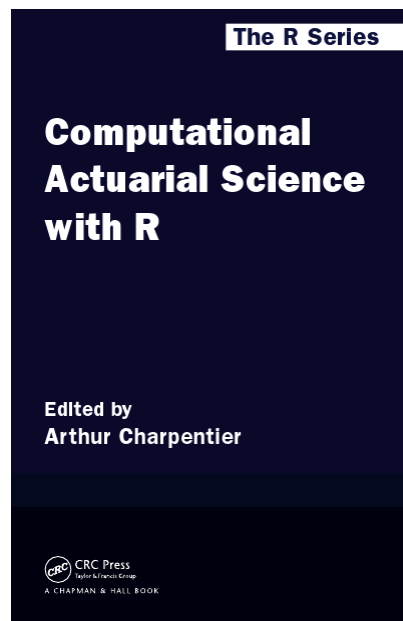
E.g. Winter temperature, in Montréal, QC



A short word on functional data



To go further...



forthcoming book entitled

Computational Actuarial Science with R