

Generalized Linear Models. Uncertainty estimation

Lecture 2c

Moving beyond typical distributions

- We know how to model
 - Normally distributed targets \rightarrow linear regression
 - Bernoulli and Multinomial targets \rightarrow logistic regression
 - What if target distribution is more complex?

Example 1: Daily Stock prices NASDAQ

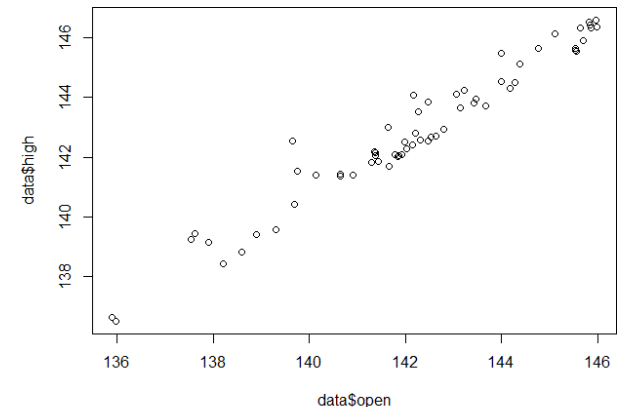
- Open
- High (within day)

Does it seem that the error is normal here?

Example 2: Number of calls to bank

- Y = Number of calls
- X = time

Endless amount of classes \rightarrow multinomial does not work... (Poisson)



Exponential family

- More advanced error distributions are sometimes needed!
- Many distributions belong to **exponential** family:
 - Normal, Exponential, Gamma, Beta, Chi-squared..
 - Bernoulli, Multinoulli, Poisson...

$$p(\mathbf{x}|\boldsymbol{\eta}) = h(\mathbf{x})g(\boldsymbol{\eta})e^{(\boldsymbol{\eta}^T u(\mathbf{x}))}$$

- Easy to find MLE and MAP
- Non-exponential family distributions: uniform, Student t

Example: Bernoulli

Generalized linear models

- Assume Y from the exponential family
- **Model** is $Y \sim EF(\mu, \dots)$, $f(\mu) = \mathbf{w}^T \mathbf{x}$
 - Alt $\mu = f^{-1}(\mathbf{w}^T \mathbf{x})$
 - f^{-1} is activation function
 - f is link function (in principle, arbitrary)
- Arbitrary f will lead to (s – dispersion parameter)

$$p(y|w, s) = h(y, s)g(\mathbf{w}, \mathbf{x})e^{\frac{b(\mathbf{w}, \mathbf{x})y}{s}}$$

- If f is a canonical link, then

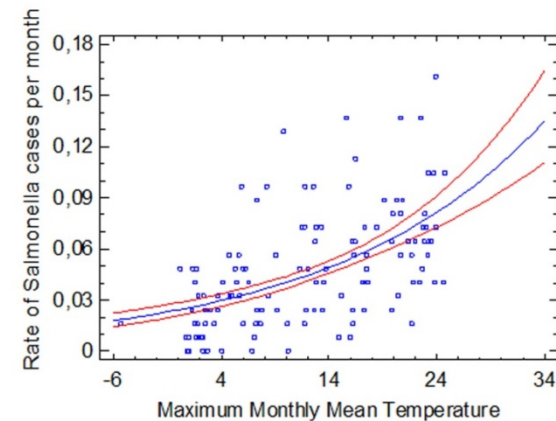
$$p(y|w, s) = h(y, s)g(\mathbf{w}, \mathbf{x})e^{\frac{(\mathbf{w}^T \mathbf{x})y}{s}}$$

Generalized linear models

- Canonical links are normally used
 - MLE computations simplify
 - MLE $\hat{w} = F(X^T Y) \rightarrow$ computations do not depend on all data but rather a summary (sufficient statistics) \rightarrow computations speed up

Example: Poisson regression

$$f^{-1}(\mu) = e^{\mu}, Y \sim \text{Poisson}(e^{w^T x})$$



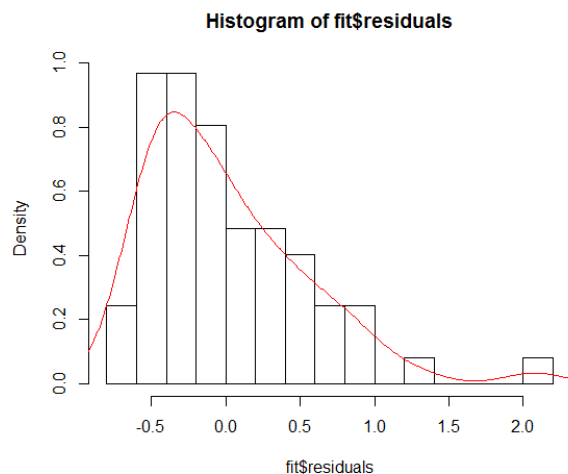
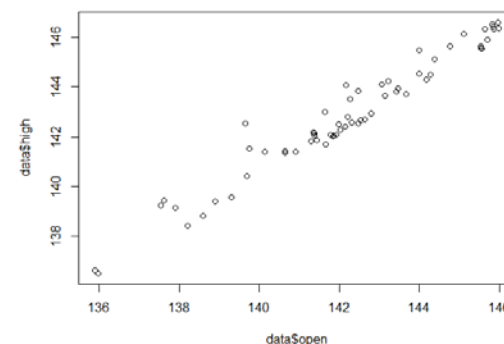
Generalized linear model: software

- Use `glm(formula, family, data)` in R

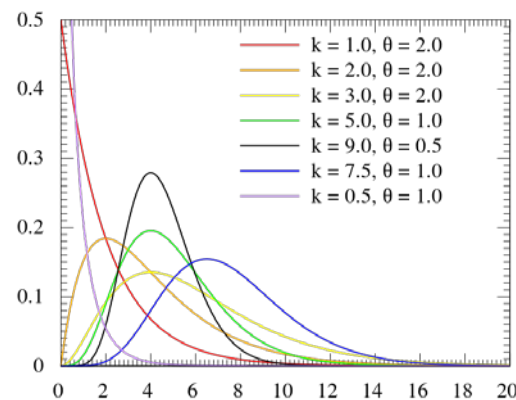
Example: Daily Stock prices NASDAQ

- Open
- High (within day)

1. Try to fit usual linear regression, study histogram of residuals

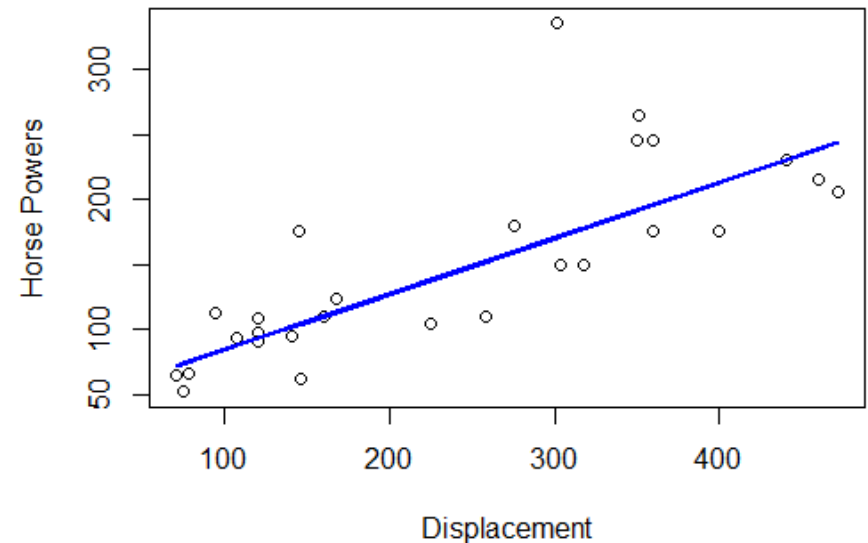


Gamma distribution: Wikipedia



Least absolute deviation regression

- Model $Y \sim \text{Laplace}(w^T X, b)$
 - Member of exponential family
- Equivalent to minimizing sum of absolute deviations
- Properties
 - Robust to outliers
 - Sensitive to changes in data
 - Multiple solutions possible
- R: package **L1pack**



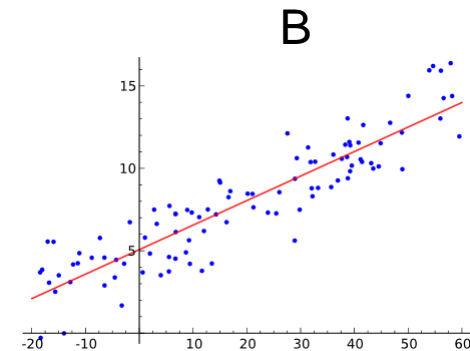
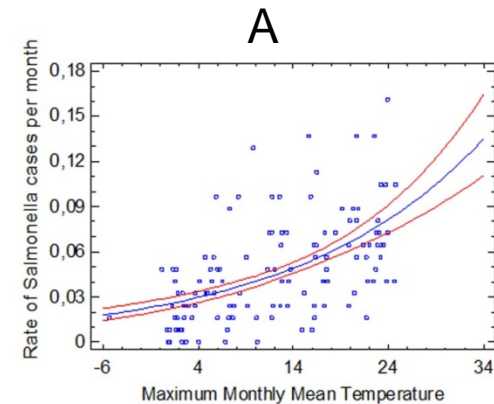
Probabilistic models

- Why it is beneficial to assume a **probabilistic** model?
- A common approach to modelling in CS and engineering:
$$y = f(x, w)$$
- f is known, w is unknown
- Fit model to data with least squares, optimization or ad hoc → find w

Probabilistic models

Arguments against deterministic models:

- The model does not really describe actual data (error is not explained)
 - No difference between modelling data A (Poisson) and B (Normal)
 - Estimation strategy for A is not good for B
- The model typically gives a **deterministic answer**, no information about uncertainty
 - "...The exchange rate tomorrow will be 8.22 ..." 😬



Probabilistic models

Probabilistic model

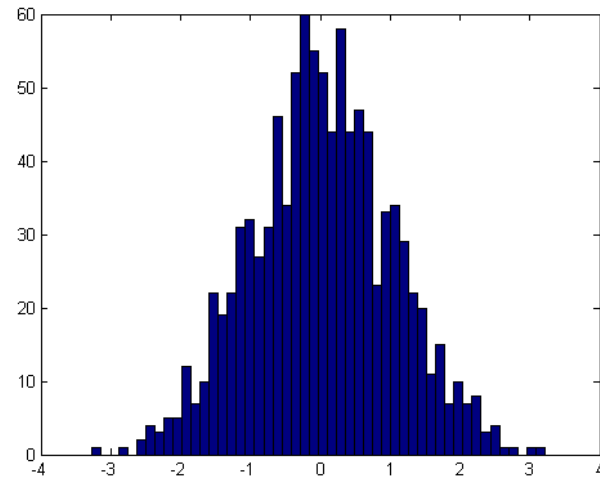
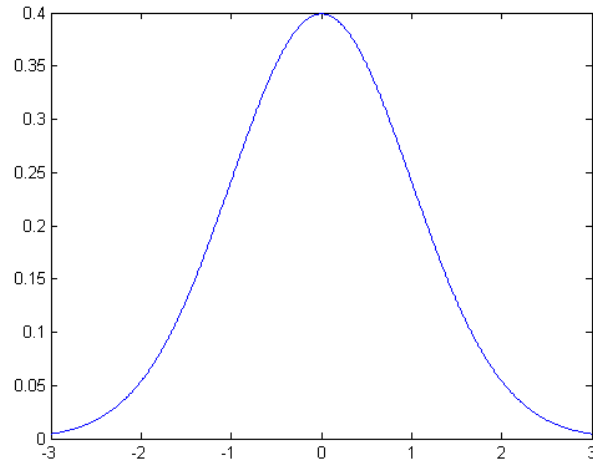
$$Y \sim \text{Distribution}(f(x, w), \theta)$$

- Data is fully explained (error as well)
- Automatic principle for finding parameters: MLE , MAP or Bayes theorem
- Automatic principle for finding uncertainty (conf. limits)
 - **Bootstrap**
 - Posterior probability
- Possibility to generate new data of the same type
 - Further testing of the model

Uncertainty estimation

- Given estimator $\hat{f} = \hat{f}(x, D)$ (or $\hat{\alpha} = \delta(D)$), how to estimate the uncertainty?
- **Answer 1:** if the distribution for data D is given, compute analytically the distribution for the estimator → derive confidence limits
 - Often difficult
 - **Example:** In simple linear regression, $\hat{\alpha}$ follows t distribution
- **Answer 2:** Use **bootstrap**

The bootstrap: general principle



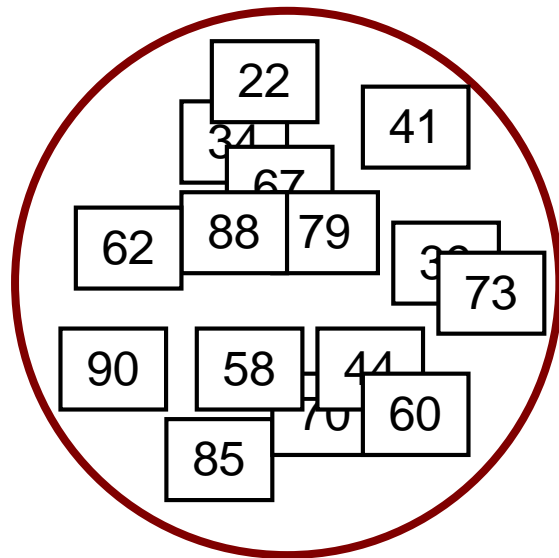
We want to determine uncertainty of $\hat{f}(D, X)$

1. Generate many different D_i from their distribution
2. Use histogram of $\hat{f}(D_i, X)$ to determine confidence limits → unfortunately can not be done (*distr of D is often unknown*)

Instead: Generate many different D_i^* from the empirical distribution (histogram)

Nonparametric bootstrap

Observed data

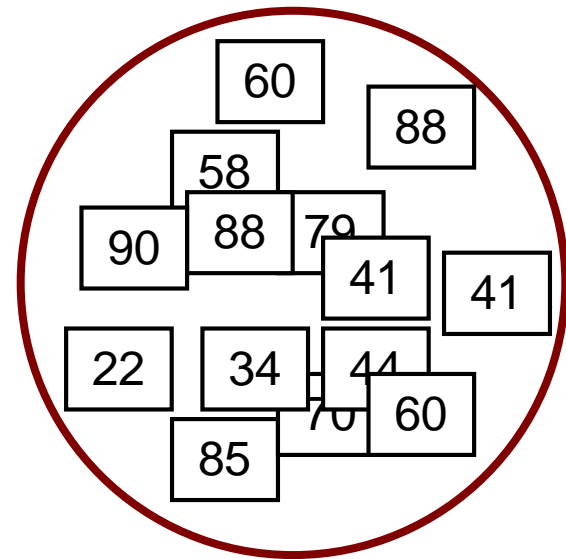


\bar{x}

Sampling with
replacement



Resampled data



$\bar{x}_1^*, \bar{x}_2^*, \dots, \bar{x}_N^*$

Nonparametric bootstrap

Given estimator $\hat{w} = \hat{f}(D)$

Assume $X \sim F(X, w)$, F and w are unknown

1. Estimate \hat{w} from data $\mathbf{D}=(X_1, \dots, X_n)$
2. Generate $\mathbf{D}_1=(X_1^*, \dots, X_n^*)$ by sampling with replacement
3. Repeat step 2 B times
4. The distribution of w is given by $\hat{f}(D_1), \dots, \hat{f}(D_B)$

Nonparametric bootstrap can be applied to any deterministic estimator, distribution-free

Parametric bootstrap

Given estimator $\hat{w} = \hat{f}(D)$

Assume $X \sim F(X, w)$, F is known and w is unknown

1. Estimate \hat{w} from data $\mathbf{D}=(X_1, \dots, X_n)$
2. Generate $\mathbf{D}_1=(X_1^*, \dots, X_n^*)$ by generating from $F(X, \hat{w})$
3. Repeat step 2 B times
4. The distribution of w is given by $\hat{f}(D_1), \dots, \hat{f}(D_B)$

Parametric bootstrap is **more** precise if the distribution form is correct

Uncertainty estimation

1. Get D_1, \dots, D_B by bootstrap
 2. Use $\hat{f}(D_1), \dots, \hat{f}(D_B)$ to estimate the uncertainty
 - Bootstrap percentile
 - Bootstrap Bca
 - ...
- Bootstrap works for all distribution types
 - Can be bad accuracy for small data sets $n < 40$ (empirical is far from true)
 - Parametric bootstrap works even for small samples

Bootstrap confidence intervals

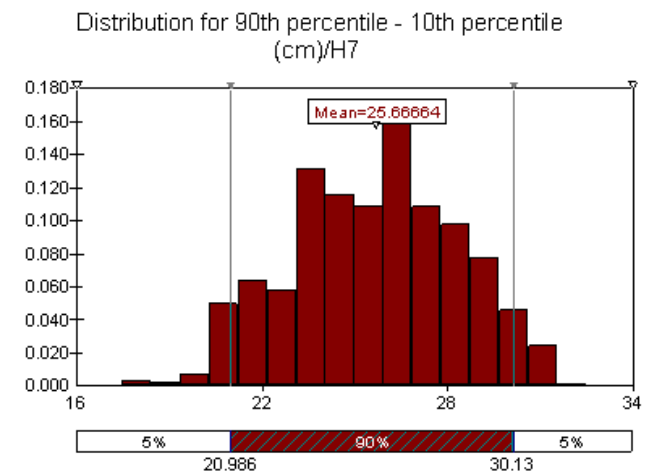
- To estimate $100(1-\alpha)$ confidence interval for w

Bootstrap percentile method

- Using bootstrap, compute $\hat{f}(D_1), \dots, \hat{f}(D_B)$, sort in ascending order, get $w_1 \dots w_B$
- Define $A_1 = \text{ceil}(B \alpha/2)$, $A_2 = \text{floor}(B - B \alpha/2)$
- Confidence interval is given by

$$(w_{A_1}, w_{A_2})$$

Look at the plot...



Bootstrap: regression context

- Model $Y \sim F(X, w)$
- Data $D = \{(Y_i, X_i), i = 1, \dots, n\}$
- Idea: produce several bootstrap sets that are similar to D

Nonparametric bootstrap:

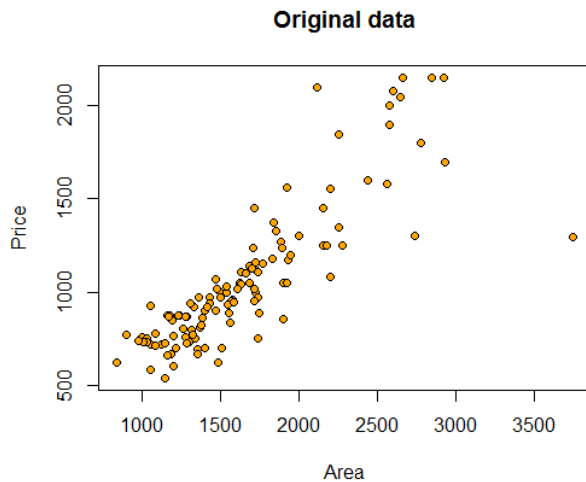
1. Using observation set D , sample **pairs** (X_i, Y_i) with replacement and get bootstrap sample D_1
2. Repeat step 1 B times \rightarrow get D_1, \dots, D_B

Uncertainty estimation

Example: Albuquerque dataset:

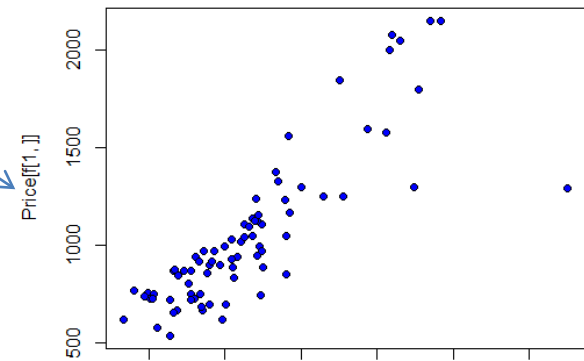
Y=Price of House

X=Area (sqft)

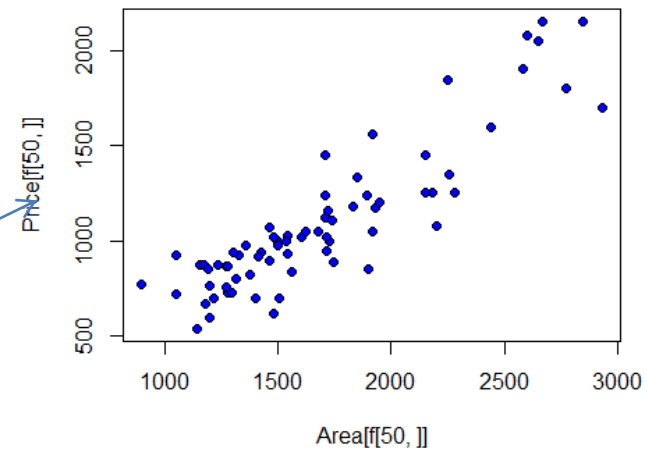


We sample data index,
from $\{1 \dots N\}$

D_1



D_{50}



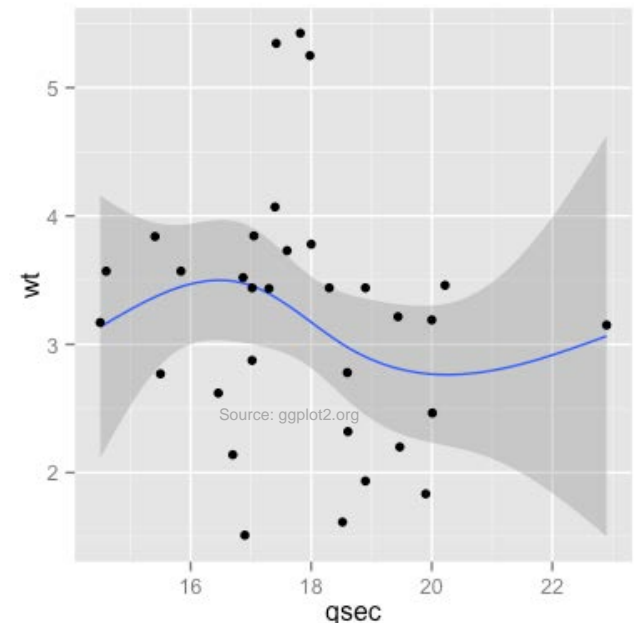
Bootstrap: regression context

Parametric bootstrap

1. Fit a model to $D \rightarrow$ get $\hat{w}(D)$.
2. Set $X_i^* = X_i$, generate $Y_i^* \sim F(X_i, \hat{w})$.
3. $D_i = \{(X_i^*, Y_i^*), i = 1, \dots, n\}$
4. Repeat step 2 B times

Confidence intervals in regression

- Given $Y \sim \text{Distribution}(y|x, w)$, $EY|X = \mu|x = f(x, w)$
 - **Example:** $Y \sim N(w^T x, \sigma^2)$, $\mu|x = f(x, w) = w^T x$
- Estimate intervals for $\mu|x = f(x, w)$ for many X , combine in a **confidence band**
- What is estimator?
 - $\mu|x = f(x, w)$



Confidence intervals in regression

Estimation

1. Compute D_1, \dots, D_B using a bootstrap
2. Fit model to $D_1, \dots, D_B \rightarrow$ estimate $\hat{w}_1, \dots, \hat{w}_B$
3. For a given X , compute $f(X, \hat{w}_1), \dots, f(X, \hat{w}_B)$ and estimate confidence interval by (percentile method)
4. Combine confidence intervals in a band

Bootstrap: R

- Package **boot**

- **Functions:**

- `boot()`
 - `boot.ci()` – 1 parameter
 - `envelope()` – many parameters

- Random random generation for parametric bootstrap:

- `Rnorm()`
 - `Runif()`
 - ...

```
boot(data, statistic, R, sim = "ordinary",  
      ran.gen = function(d, p) d, mle = NULL,...)
```

Bootstrap: R

Nonparametric bootstrap:

- Write a function *statistic* that depends on *dataframe* and *index* and returns the estimator

```
library(boot)
data2=data[order(data$Area),]#reordering data according to Area

# computing bootstrap samples
f=function(data, ind){
  data1=data[ind,]# extract bootstrap sample
  res=lm(Price~Area, data=data1) #fit linear model
  #predict values for all Area values from the original data
  priceP=predict(res,newdata=data2)
  return(priceP)
}
res=boot(data2, f, R=1000) #make bootstrap
```


Bootstrap: R

Parametric bootstrap:

- Compute value mle that estimates model parameters from the data
- Write function *ran.gen* that depends on *data* and *mle* and which generates new data
- Write function *statistic* that depend on *data* which will be generated by *ran.gen* and should return the estimator

Bootstrap

```
mle=lm(Price~Area, data=data2)

rng=function(data, mle) {
  data1=data.frame(Price=data$Price, Area=data$Area)
  n=length(data$Price)
  #generate new Price
  data1$Price=rnorm(n,predict(mle, newdata=data1),sd(mle$residuals))
  return(data1)
}

f1=function(data1){
  res=lm(Price~Area, data=data1) #fit linear model
  #predict values for all Area values from the original data
  priceP=predict(res,newdata=data2)
  return(priceP)
}

res=boot(data2, statistic=f1, R=1000, mle=mle,ran.gen=rng, sim="parametric")
```

Uncertainty estimation: R

- Bootstrap confidence bands for linear model

```
e=envelope(res) #compute confidence bands
```

```
fit=lm(Price~Area, data=data2)
```

```
priceP=predict(fit)
```

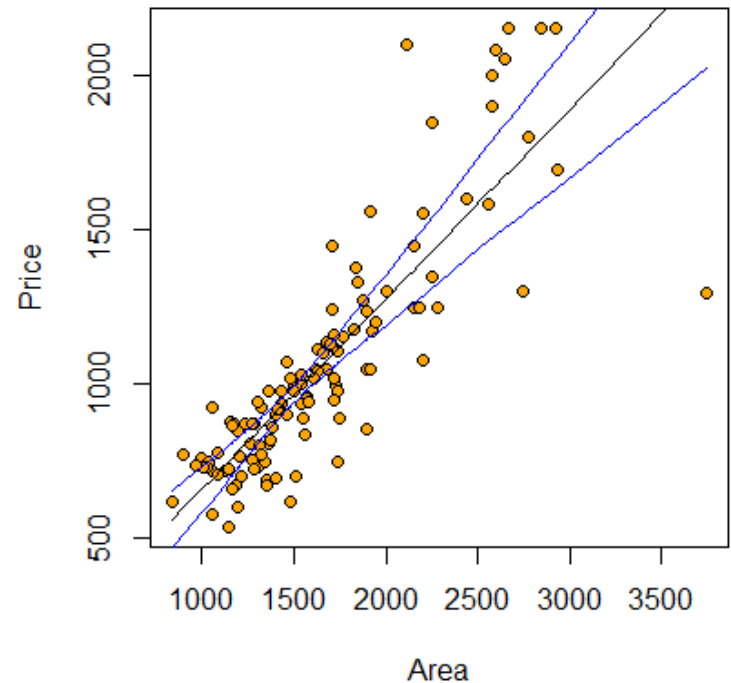
```
plot(Area, Price, pch=21, bg="orange")
```

```
points(data2$Area,priceP,type="l") #plot fitted line
```

```
#plot confidence bands
```

```
points(data2$Area,e$point[2,], type="l", col="blue")
```

```
points(data2$Area,e$point[1,], type="l", col="blue")
```



Prediction bands

- Confidence interval for $Y|X$ = interval for mean $EY|X$
- Prediction interval for $Y|X$ = interval for $Y|X$

$$Y \sim \text{Distribution}(x, w)$$

Prediction band for parametric bootstrap

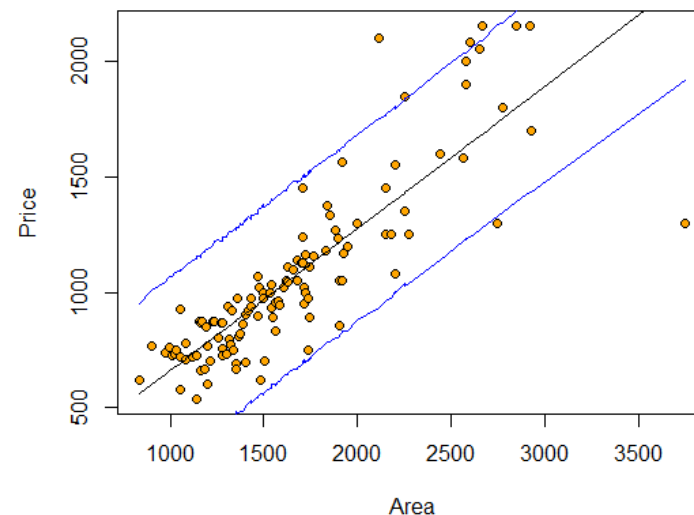
1. Run parametric bootstrap and get D_1, \dots, D_B
2. Fit the model to the data and get $\hat{w}(D_1), \dots, \hat{w}(D_B)$
3. For each X , generate from $\text{Distribution}(X, \hat{w}(D_1)), \dots, \text{Distribution}(X, \hat{w}(D_B))$ and apply percentile method
4. Connect the intervals \rightarrow get the band

Estimation of the model quality

Example: parametric bootstrap

```
mle=lm(Price~Area, data=data2)

f1=function(data1){
  res=lm(Price~Area, data=data1) #fit
  linear model
  #predict values for all Area values
  from the original data
  priceP=predict(res,newdata=data2)
  n=length(data2$Price)
  predictedP=rnorm(n,priceP,
sd(mle$residuals))
  return(predictedP)
}
res=boot(data2, statistic=f1, R=10000,
mle=mle,ran.gen=rng, sim="parametric")
```



Why wider band?