# Chapter 8

# **Nonparametric Regression**

Date $\Rightarrow$ distribution.

$\rightarrow$ assume some family, for example, $N(\mu, \sigma^2) \Rightarrow \hat{\mu}, \hat{\sigma}^2$

$\rightarrow$ without any distributional assumptions $\Rightarrow$ form from the data.

Universität Augsburg

# Nonparametric regression

$$\{(X_i, Y_i)\}, \ i = 1, \ldots, n; \quad \boldsymbol{X} \in \mathbb{R}^{J+1}, Y \in \mathbb{R}$$

- Engel curve: $X$ = net-income, $Y$ = expenditure
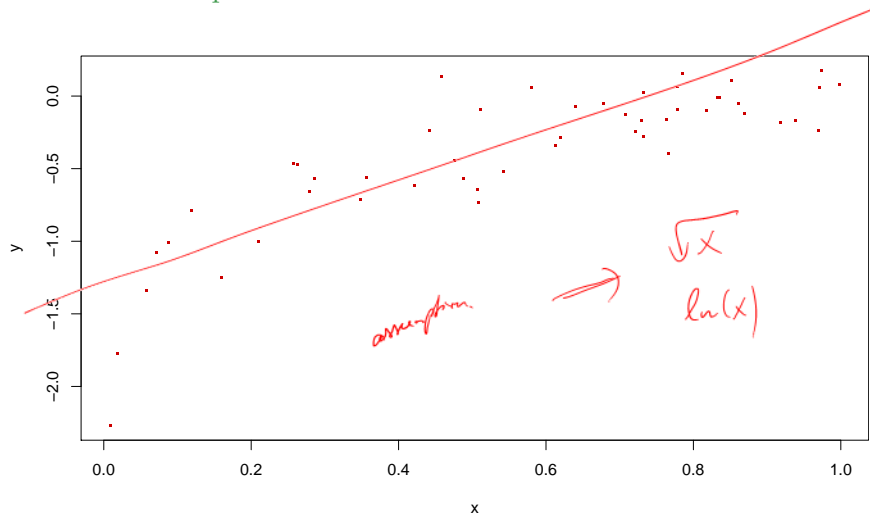
$$Y = m(X) + \varepsilon$$

- CHARN model: time series of the form

$$Y_t = m(Y_{t-1}) + \sigma(Y_{t-1})\xi_t$$

# Let us have a sample of size $n = 50$ from an unknown model.

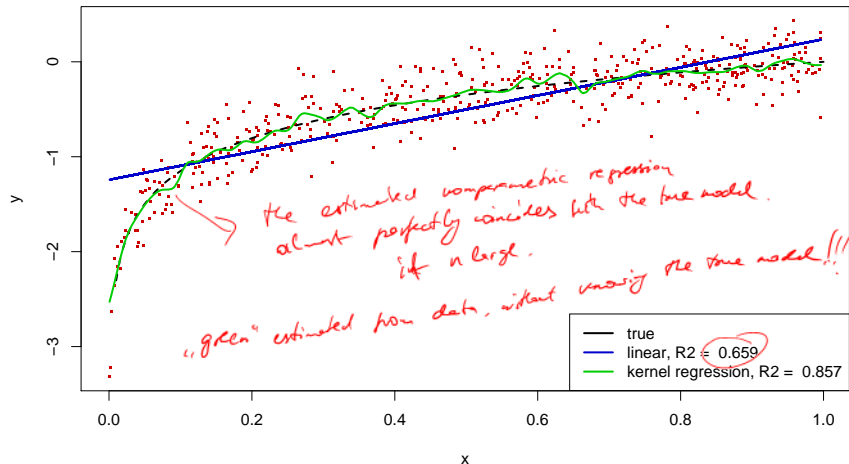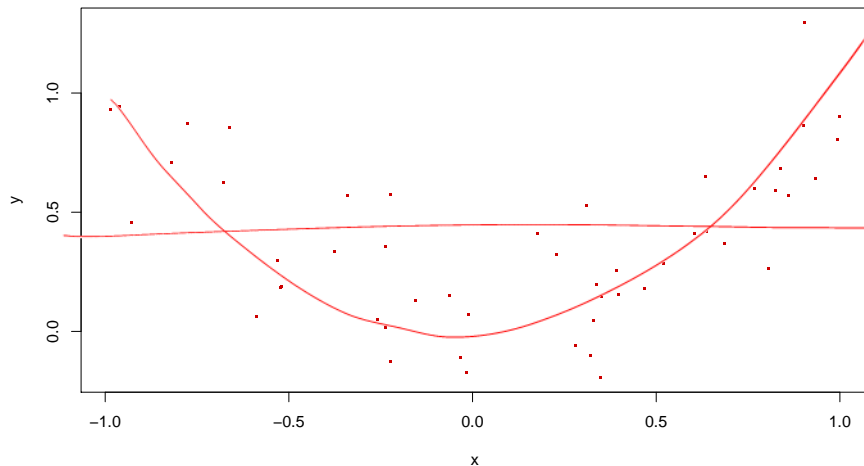# The data is simulated from the model

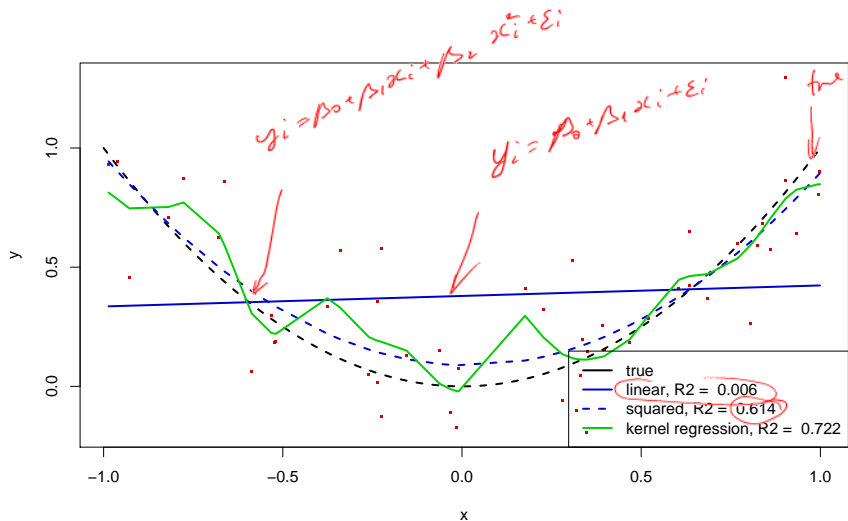$$y_i = 0.5 \ln(x_i) + 0.2 u_i, \qquad u_i \sim N(0,1).$$
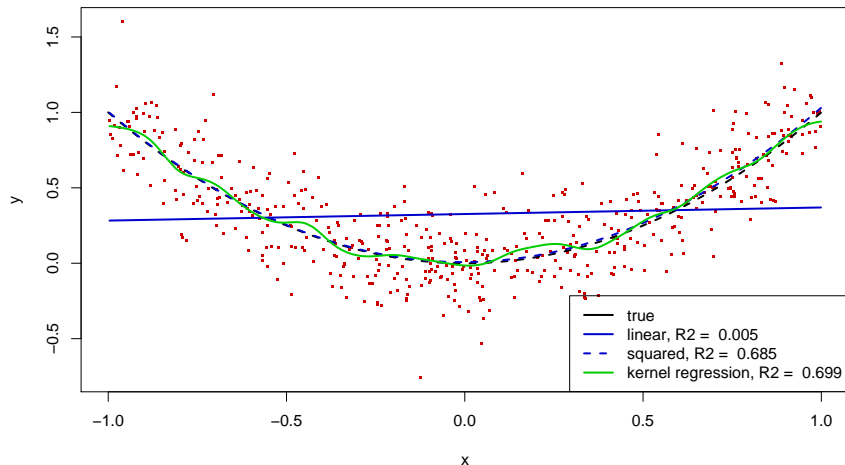
For $n = 500$

# Another example.

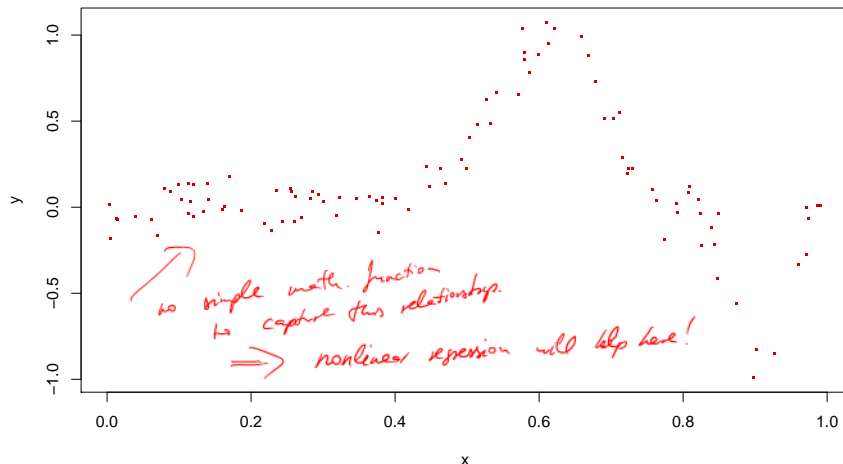# The data is simulated from the model

$$y_i = x_i^2 + 0.2u_i, \qquad u_i \sim N(0,1).$$
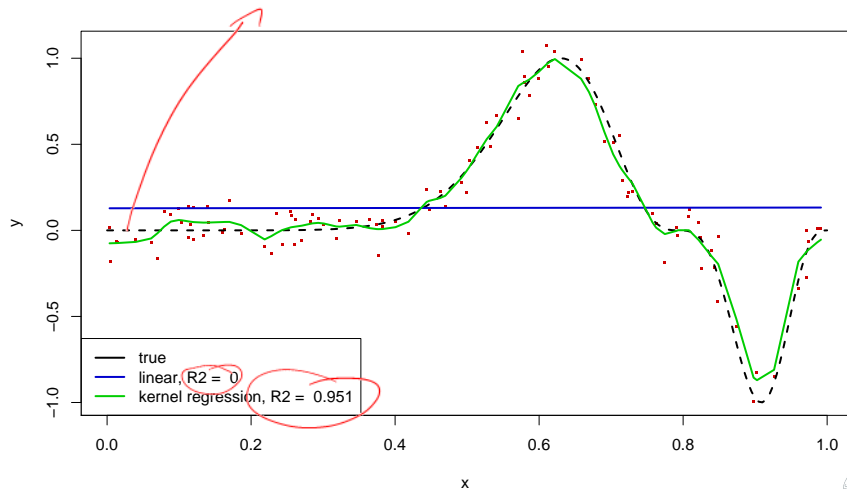
... and with $n = 500$.

# Hopeless example ...

... simulated from

$$y_i = \{sin(2\pi x_i^3)\}^3 + 0.1u_i, \qquad u_i \sim N(0,1).$$

*1st: nonparametric estimator of distribution.*

*2nd: nonparametric estimator of regression.*

# Kernel density estimator

**First:** the procedure requires a non-parametric estimator of a density.

Here: DAX30 returns, 20 years of daily data, 5217 observations with normal and $t$-densities $\leadsto$ poor fit in the middle and in the tails



*N*

*many returns, close to 0.*

$\mu = 0.038$
$\sigma = 1.367$

*estimator of the distribution obtained directly from data.*

*t.*

*large distance between N & t and the histogram.*

*⟹ N & t are not appropriate.*

*large losses*

*⟨ 0*

*very large return ≥ 0*

*large profits*

*the same probabilities, but different x-values*

## Drawbacks of the histogram

- constant over intervals, step function
- results depend strongly on origin
- binwidth choice ⟹ *width of rectangles?*
- (slow rate of convergence)

*starting point*

Universität
Augsburg

# ... and with a kernel density estimator



*is light to reflect way returns which are close to 0.*

$\mu$ = 0.038
$\sigma$ = 1.367

Legend:
- Normal density
- Kernel density
- t–density, 4df

*Better fit expand to N art.*

# Kernel Density Estimation

**KDE as a generalization of the histogram**

Idea of the histogram:

$$\frac{1}{n \cdot \text{interval length}} \#\{\text{obs. that fall into a small interval CONTAINING } x\}$$

Idea of the kernel density:

$$\frac{1}{n \cdot \text{interval length}} \#\{\text{obs. that fall into a small interval AROUND } x\}$$

*indicator function.*

$$
\begin{aligned}
\widehat{f}_h(x) &= \frac{1}{2hn} \sum_{i=1}^{n} I\left(\left|\frac{x - X_i}{h}\right| \le 1\right) \\
&= \frac{1}{2hn} \#\{X_i \text{ in interval around } x\}
\end{aligned}
$$

Universität Augsburg

kernel density estimate (KDE)

$$\widehat{f}_h(x) = \frac{1}{nh} \sum_{i=1}^{n} K\left(\frac{x - X_i}{h}\right)$$

*bell-shaped*

with kernel function $K(u)$
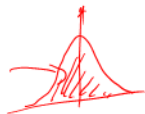
Required properties of kernels

*the square under $K = 1$*

- $K(\bullet)$ is a density function:

$$\int_{-\infty}^{\infty} K(u)du = 1 \quad \text{and} \quad K(u) \geq 0$$

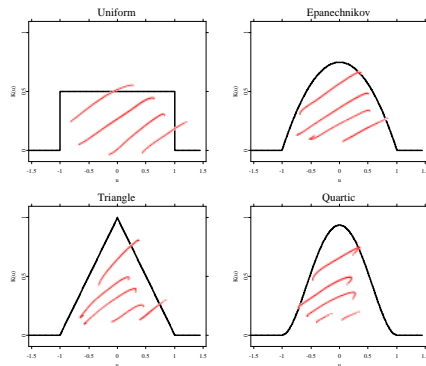- $K(\bullet)$ is symmetric:

$$\int_{-\infty}^{\infty} uK(u)du = 0$$

# Different Kernel Functions

| Kernel | $K(u)$ |
|---:|:---|
| Uniform | $\frac{1}{2}I(|u| \le 1)$ |
| Triangle | $(1 - |u|)I(|u| \le 1)$ |
| Epanechnikov | $\frac{3}{4}(1 - u^2)I(|u| \le 1)$ |
| Quartic | $\frac{15}{16}(1 - u^2)^2 I(|u| \le 1)$ |
| Triweight | $\frac{35}{32}(1 - u^2)^3 I(|u| \le 1)$ |
| Gaussian | $\frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}u^2)$ |
| Cosinus | $\frac{\pi}{4} \cos(\frac{\pi}{2}u)I(|u| \le 1)$ |

*(handwritten annotations: "theoretical properties." pointing to Epanechnikov row; "because of the round." pointing to Gaussian row)*

Table: Kernel functions

Figure: Some kernel functions: Uniform (top left), Triangle (bottom left), Epanechnikov (top right), Quartic (bottom right)

**Example:** Construction of the KDE

consider the KDE using a Gaussian kernel

$$\widehat{f}_h(x) = \frac{1}{nh} \sum_{i=1}^{n} K\left(\frac{x - X_i}{h}\right)$$

$$= \frac{1}{nh} \sum_{i=1}^{n} \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}u^2)$$

here we have

*the point where the density must be estimated*

$$u = \frac{x - X_i}{h}$$

*data point (sample)*

*bandwidth*

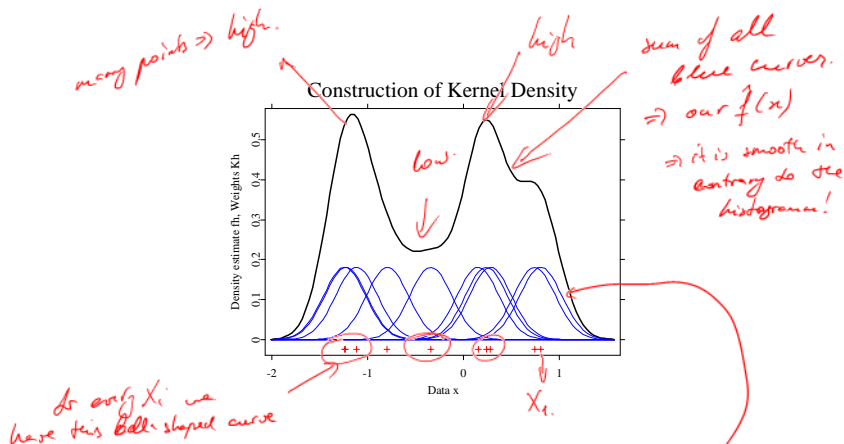Universität Augsburg

*handwritten annotations:*

many points ⇒ high.

high

sum of all blue curves.
⇒ our $f(x)$
⇒ it is smooth in contrary to the histogram!

low.

Construction of Kernel Density

for every $X_i$ we have this bell-shaped curve

$X_1$.

$X_1 \Rightarrow \dfrac{1}{nh\sqrt{2\pi}} exp\left(\dfrac{x - X_1}{h}\right)^2 \Rightarrow$

Figure: Kernel density estimate as a sum of bumps

## Plotting kernel estimators in R

*band-idth* ↓

*"method"* ↓

*R functions.* ↓

```
density(x, bw = "nrd0", adjust = 1,
        kernel = c("gaussian", "epanechnikov", "rectangular",
                   "triangular", "biweight", "cosine", "optcosine"),
        weights = NULL, window = kernel, width, give.Rkern = FALSE,
        n = 512, from, to, cut = 3, na.rm = FALSE, ...)


Value:
x: the 'n' coordinates of the points where the density is
            estimated.
y: the estimated density values.  These will be non-negative,
            but can be zero.
bw: the bandwidth used.
```

```
> x = rnorm(100)                    100 obs from N(0,1)
> k = density(x)
> k
Data: x (100 obs.);     Bandwidth 'bw' = 0.3029
       x                 y
 Min.   :-3.4975   Min.   :0.000175
 1st Qu.:-1.7048   1st Qu.:0.023764
 Median : 0.0879   Median :0.061481
 Mean   : 0.0879   Mean   :0.139314
 3rd Qu.: 1.8806   3rd Qu.:0.248836
 Max.   : 3.6733   Max.   :0.443788                 grid on the x-axis
> k$x
  [1] -3.497498092 -3.483465207 -3.469432322 -3.455399437 -3.441366553
  [6] -3.427333668 -3.413300783 -3.399267898 -3.385235013 -3.371202128
  ...
> k$y
  [1] 0.0002050555 0.0002360509 0.0002706135 0.0003111481 0.0003565894
  [6] 0.0004070709 0.0004629191 0.0005277777 0.0005996858 0.0006789000
  ...                                    ↑ f̂(grid).
```
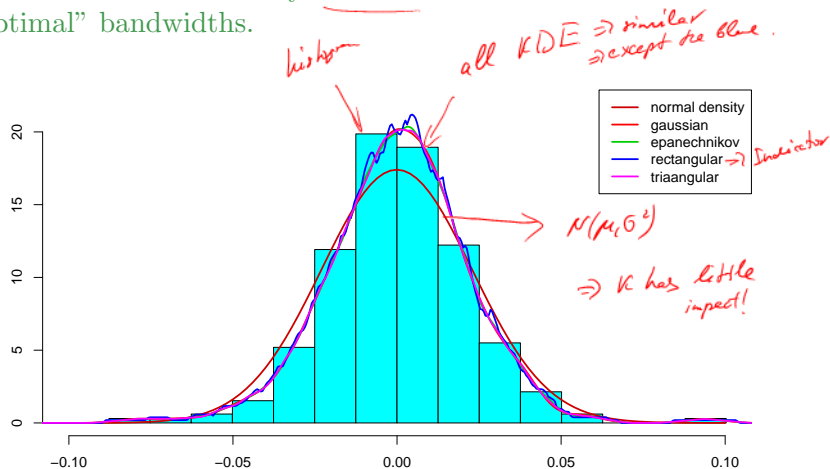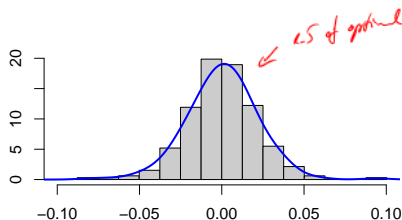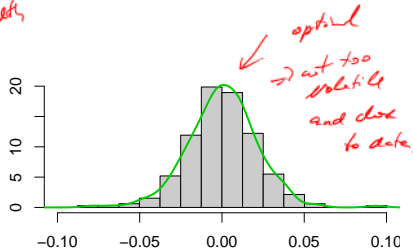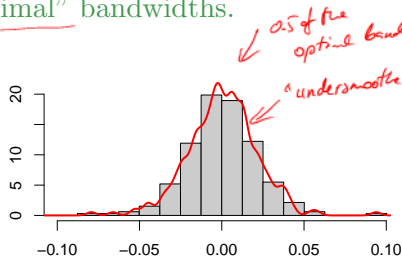
Kernel estimator of the density of DJ30 returns with different kernels and "optimal" bandwidths.
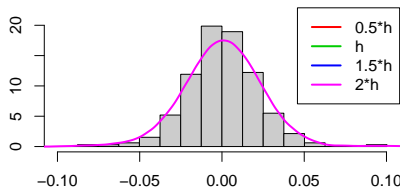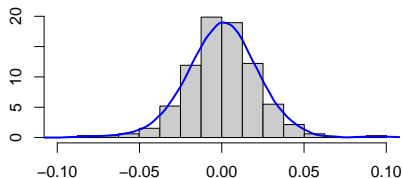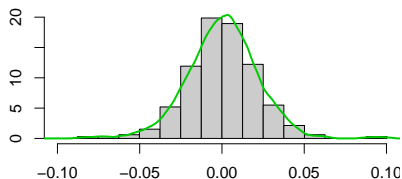
```
z1=density(rdj, bw="nrd0",  kernel="epanechnikov")

pdf("ch2_hist_dax.pdf", width=9, height=5, onefile=FALSE);
truehist(rdj,  prob=TRUE, xlab="",xlim=c(-0.07,0.07), col="grey80");
matplot(z1$x,z1$y, add=T, type="l", lty=1, col=2, lwd=2);
dev.off()
```

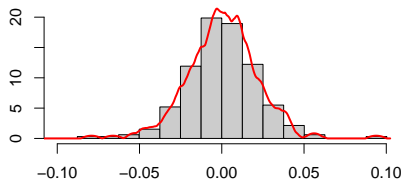Gaussian kernel for the density of DJ30 returns with different adjusted "optimal" bandwidths.

Epanechnikov kernel for the density of DJ30 returns with different adjusted "optimal" bandwidths.

*Almost no visible changes if we move from Gaussian to Epan*

# Epanechnikov kernel for the density of DAX30 returns with the bandwidths 0.01, 0.001, 0.0001.

Universität
Augsburg

# Example: time-variation of the distribution of DAX returns (Epanechnikov kernel)

# Example: time-variation of the distribution of DAX returns (normal density)

European Community Household Panel for the period 1994-2001. Data for Germany, net household income, ca. 48000 observations and Epanechnikov kernel.

European Community Household Panel for the period 1994-2001. Data for Germany, net household income, ca. 48000 observations and gaussian kernel.

Universität
Augsburg

European Community Household Panel for the period 1994-2001. Data for Germany, net household income, ca. 48000 observations and gaussian density .

**(Asymptotic) statistical properties of KDE** bias of the kernel density estimator

$E(\hat{\Theta}) \quad - \Theta \quad = 0 \quad$ if unbiased

$$
\begin{aligned}
Bias\left\{\widehat{f}_h(x)\right\} &= E\left[\widehat{f}_h(x)\right] - f(x) \\
&= E\left[\frac{1}{nh}\sum_{i=1}^{n} K\left(\frac{x-X_i}{h}\right)\right] - f(x) \\
&\approx \frac{h^2}{2}f''(x)\mu_2(K) \quad \text{for } h \to 0
\end{aligned}
$$

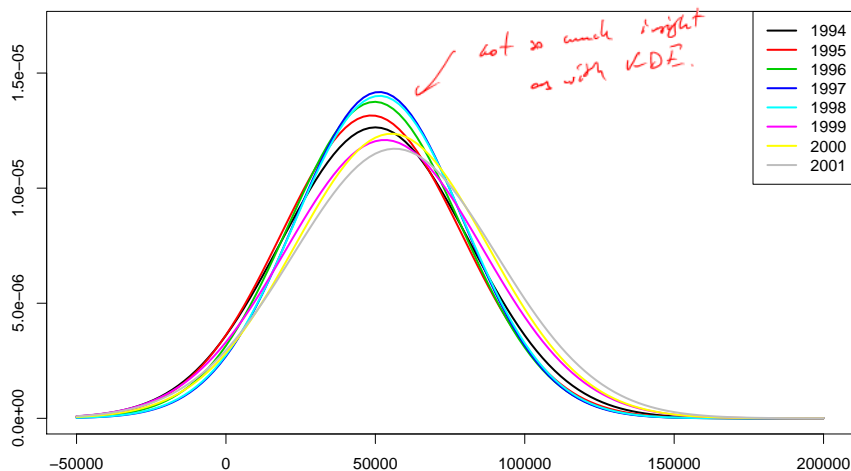Bias disappears if bandwidth is getting smaller!!!

characteristic of $K$
characteristic of the true density

where $\mu_2(K) = \int_{-\infty}^{\infty} u^2 K(u)du$.

$$
\begin{aligned}
Var\left[\widehat{f}_h(x)\right] &= Var\left[\frac{1}{nh}\sum_{i=1}^{n} K\left(\frac{x-X_i}{h}\right)\right] \\
&\approx \frac{1}{nh}\|K\|_2^2 f(x) \quad \text{for } nh \to \infty
\end{aligned}
$$

true density

where $\|K\|_2^2 = \int_{-\infty}^{\infty} \{K(u)\}^2 du$.

characteristic of $K$

sample size. $\to$ $h$ increases $\Rightarrow$ Var

trade off: $h \nearrow$ $\Rightarrow$ Var $\searrow$ if $h \searrow$ $\Rightarrow$ Bias $\searrow$ $\}$ $h_{opt}$ ?

Universität Augsburg

## How to choose the bandwidth for the KDE?

find the bandwidth which minimizes the $MISE$

$$
MISE\left\{\widehat{f}_h(x)\right\} = E\left[\int_{-\infty}^{\infty}\{\widehat{f}_h(x) - f(x)\}^2 dx\right]
$$

mean integrated
squared error

$$
= \int_{-\infty}^{\infty} MSE[\widehat{f}_h(x)]dx \qquad \text{known / fixed}
$$

$$
\approx \frac{1}{n\,h}\|K\|_2^2 + \frac{h^4}{4}\mu_2(K)^2\|f''\|_2^2 = AMISE\left\{\widehat{f}_h(x)\right\}
$$

$\longrightarrow$ min w.r.t. $h$

and thus

$$
h_{opt} = \left(\frac{\|K\|_2^2}{\|f''\|_2^2\mu_2(K)^2 n}\right)^{1/5} \sim n^{-1/5}
$$

are determined by $K$

Problematic: to estimate $f$ you need $h$, but to compute $h$ you need $f$ !

$\|f''\| \Longrightarrow$ norm. density $\Longrightarrow$ Silverman's rule of thumb

Universität Augsburg

## Multivariate KDE

$d$-dimensional data and $d$-dimensional kernel

$$\boldsymbol{X} = (X_1, \ldots, X_d)^\top, \quad \mathcal{K} : \mathbb{R}^d \to \mathbb{R}_+$$

- multivariate kernel density estimator (simple)

$$\widehat{f}_h(\boldsymbol{x}) = \frac{1}{n} \sum_{i=1}^{n} \frac{1}{h^d} \mathcal{K}\left(\frac{x_i - X_i}{h}\right)$$

each component is scaled equally.

- multivariate kernel density estimator (more general)

*multivariate kernel $\mathcal{K}$.*

$$\widehat{f}_h(\boldsymbol{x}) = \frac{1}{n} \sum_{i=1}^{n} \frac{1}{h_1 \cdot \ldots \cdot h_d} \mathcal{K}\left(\frac{x_1 - X_1}{h_1}, \ldots, \frac{x_d - X_d}{h_d}\right)$$

*individual bandwidths for each dimension.*

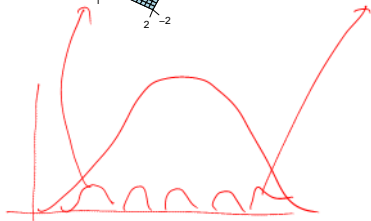- multivariate kernel density estimator (most general)

$$\widehat{f}_{\mathbf{H}}(\boldsymbol{x}) = \frac{1}{n} \sum_{i=1}^{n} \frac{1}{\det(\mathbf{H})} \mathcal{K} \left\{ \mathbf{H}^{-1}(\boldsymbol{x} - \boldsymbol{X}_i) \right\}$$

where $\mathbf{H}$ is a (symmetric) bandwidth matrix.

Each component is scaled separately, correlation between components can be handled.

**Example:** product Epanechnikov kernel with bandwidth matrices

$$\mathbf{H} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \mathbf{H} = \begin{pmatrix} 0.5 & 0 \\ 0 & 1 \end{pmatrix}, \mathbf{H} = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix}$$

Universität Augsburg

**Example:** radially symmetric Epanechnikov kernel with bandwidth matrices $\mathbf{H} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$, $\mathbf{H} = \begin{pmatrix} 0.5 & 0 \\ 0 & 1 \end{pmatrix}$, $\mathbf{H} = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix}$

**Example:** bandwidth matrix $\mathbf{H} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$

*more weights on the angles of the distribution*



Figure: Contours from bivariate product (left) and bivariate radially symmetric (right) Epanechnikov kernel

**Example:** bandwidth matrix $\mathbf{H} = \begin{pmatrix} 1 & 0 \\ 0 & 0.5 \end{pmatrix}$



Figure: Contours from bivariate product (left) and bivariate radially symmetric (right) Epanechnikov kernel

**Example:** bandwidth matrix $\mathbf{H} = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix}^{1/2}$
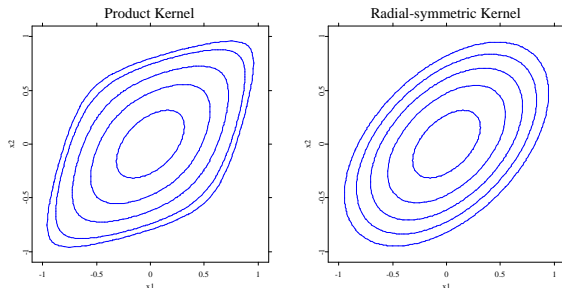


Figure: Contours from bivariate product (left) and bivariate radially symmetric (right) Epanechnikov kernel

# Kernel properties

$\mathcal{K}$ is a density function

$$\int_{\mathbb{R}^d} \mathcal{K}(\boldsymbol{u})d\boldsymbol{u} = 1, \ \mathcal{K}(\boldsymbol{u}) \geq 0$$

$\mathcal{K}$ is symmetric

$$\int_{\mathbb{R}^d} \boldsymbol{u}\mathcal{K}(\boldsymbol{u})d\boldsymbol{u} = 0_d$$

Italian GDP growth panel for 21 regions covering the period 1951-1998 (millions of Lire, 1990=base). There are 1008 observations in total.

```
data("Italy")
fhat = npcdens(gdp ~ year, tol=0.1, ftol=0.1, data=Italy)
summary(fhat)
plot(fhat, view="fixed", main="", theta=300, phi=50)
```

⟹ 2-dih. data

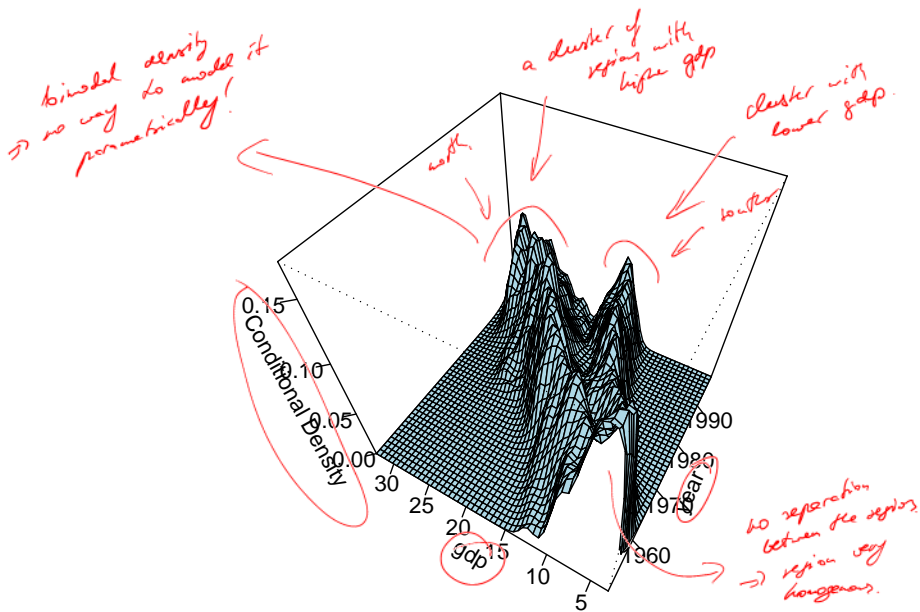```
Conditional Density Data: 1008 training points, in 2 variable(s)
(1 dependent variable(s), and 1 explanatory variable(s))
                           gdp
Dep. Var. Bandwidth(s): 0.697371
                          year
Exp. Var. Bandwidth(s): 0.6725248

Bandwidth Type: Fixed
Log Likelihood: -2550.493
```

Waiting time between eruptions and the duration of the eruption for the Old Faithful geyser in Yellowstone National Park, Wyoming, USA.

```
library("datasets");library("np");
data("faithful"); attach(faithful);

bw = npudensbw(dat=faithful)
summary(bw)

npplot(bws=bw, xtrim=-0.2)
detach("faithful")

Data (272 observations, 2 variable(s)):
Bandwidth Selection Method: Maximum Likelihood Cross-Validation
Bandwidth Type: Fixed
Objective Function Value: 4.197484 (achieved on multistart 1)

Var. Name: eruptions Bandwidth: 0.1470088 Scale Factor: 0.327852
Var. Name: waiting   Bandwidth: 2.925438 Scale Factor: 0.5477395
```

**[theta= 0, phi= 10]**

Universität Augsburg

## Summary for KDE

- The KDE does not depend on the starting points of the classes.
- The resulting estimator is a smooth and continuous function.
- The estimator heavily depends on the bandwidth. $\Rightarrow h_{opt}$
- The estimator is robust to different choices of the kernel function.
- The estimator is biased in general.   $\rightarrow$ Gauss, Epane
- Decreasing bandwidth implies smaller bias, but larger variance.

  $\rightarrow$ two slides with formulas.

  $h_{opt} \Rightarrow$ trade-off ( as for usual parameters )

Universität
Augsburg

*LR: $Y_i = \beta_0 + \beta_1 X_{ij} + \varepsilon_i \Rightarrow$ we approx. m by a linear function.*

*$E(l_i'|X_i) = \beta_0 + \beta_1 X_i$*

# Univariate nonparametric regression

*Here: estimate $m(\cdot)$ from data, without fixing it.*

model

*on analogy $E(Y_i|X_i = x_i)$*

$$Y_i = m(X_i) + \varepsilon_i, \quad i = 1, \ldots, n$$

$m(\bullet)$ smooth regression function, $\varepsilon_i$ i.i.d. error terms with $E\varepsilon_i = 0$

we aim to estimate the conditional expectation of $Y$ given $X = x$

*cond. expect — leads to a conditional density function. joint density bivariate.*

$$m(x) = E(Y|X = x) = \int y \, f(y|x) \, dy = \int y \, \frac{f(x,y)}{f_X(x)} \, dy$$

*marginal density of X*

where $f(x, y)$ denotes the joint density of $(X, Y)$ and $f_X(x)$ the marginal density of $X$

*$= \dfrac{\left(\int y \, \hat{f}(x,y) dy\right)}{\hat{f}_X(x)}$*

*$\hat{f}$ and $\hat{f}_X$ can be estimated by KDE*

*$\Rightarrow \hat{m}$ without any distributional f functional assumptions !!!*

*$E(Y) = \int y f(y) dy$ Def — density of Y*

*$P(A|B) = \dfrac{P(A \cap B)}{P(B)}$ Def of conditional probability*

## Nadaraya-Watson Estimator

idea: $(X_i, Y_i)$ have a joint pdf, so we can estimate $m(\bullet)$ by a multivariate kernel estimator

*KDE.*

*estimator for f(x,y)*

$$\widehat{f}_{h,\widetilde{h}}(x,y) = \frac{1}{n}\sum_{i=1}^{n}\frac{1}{h}K\left(\frac{x-X_i}{h}\right)\frac{1}{\widetilde{h}}K\left(\frac{y-Y_i}{\widetilde{h}}\right)$$

*part for X*   *for y.*

and therefore $\underbrace{\int y\,\widehat{f}_{h,\widetilde{h}}(x,y)dy} = \frac{1}{n}\sum_{i=1}^{n}K_h(x-X_i)Y_i$ ⟹ *what is left from by integral.*

resulting estimator:

*estimated regression function*

$$\widehat{m}_h(x) = \frac{n^{-1}\sum_{i=1}^{n}K_h(x-X_i)Y_i}{n^{-1}\sum_{i=1}^{n}K_h(x-X_i)} = \frac{\widehat{r}_h(x)}{\widehat{f}_h(x)}$$

*$\hat{f}(x,y)$*

*no function assumptions !!!*

*$K_h(u)=\frac{1}{h}K(\frac{u}{h})$*

*$\hat{f}_X(x) \rightarrow KDE$*

*green points on the next slide.*

## Observations and weights $W_{hi}$ for $h = 0.1$, Gaussian kernel and different values of $x = x_0$.

Example: happiness

The happiness of nations (measured as average happiness of the citizens) depends of the average income per capita.

Data: 62 nations

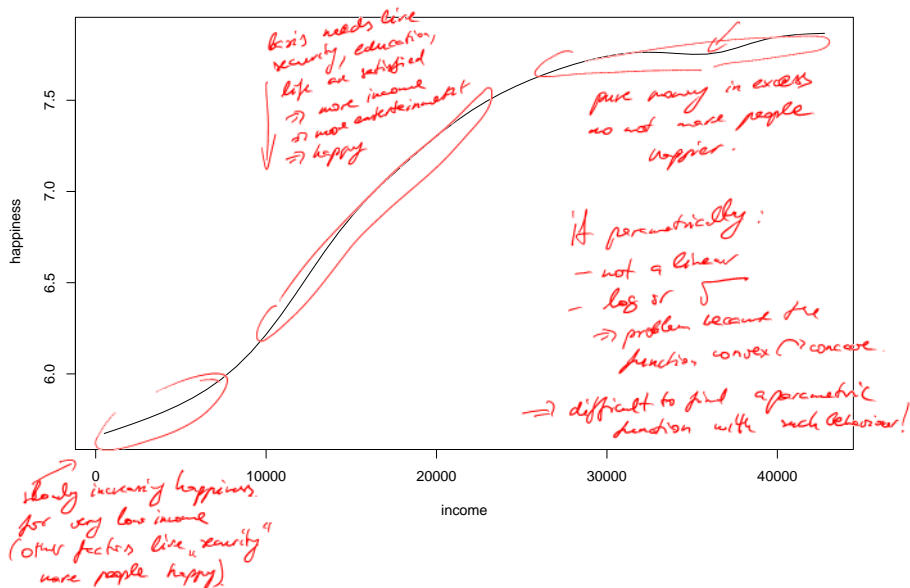$$\text{happiness} = m(\text{income}) + \varepsilon.$$

## Nonparametric regression in R

```
library("np")
all = read.table("ch3_happiness.txt");
happiness = all[,1];
income = all[,2];

bw = npregbw(formula=happiness ~ income, lt="lc")
model = npreg(bws=bw);

npplot(bw, type="l");
```

*← one function to determine the bandwidth*

*→ one function to compute the estimator.*

Universität
Augsburg

## Example: wage in Canada

Canadian cross-section wage data consisting of a random sample taken from the 1971 Canadian Census Public Use Tapes for male individuals having common education (Grade 13). There are n = 205 observations in total, and 2 variables, the logarithm of the individual's wage (logwage) and their age (age).

*to guarantee*
*> 0 forecasts*
*and more*
*residuals*
*symmetric*

$$log(\text{wage}) = m(\text{age}) + \varepsilon$$

$$logwage = \beta_0 + \beta_1 \cdot wage + \varepsilon_i$$

```
data("cps71")
model.par = lm(logwage ~ age , data = cps71)
summary(model.par)

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 13.021902   0.144705   89.99  < 2e-16 ***
age          0.012046   0.003554    3.39  0.00084 ***

Residual standard error: 0.6206 on 203 degrees of freedom
Multiple R-squared: 0.05357,    Adjusted R-squared: 0.04891
F-statistic: 11.49 on 1 and 203 DF,  p-value: 0.0008407
```

$$logwage = \beta_0 + \beta_1 \cdot age + \beta_2 \cdot age^2 + \varepsilon_i$$

```
model.par2 = lm(logwage ~ age+I(age^2), data = cps71)
summary(model.par2
```
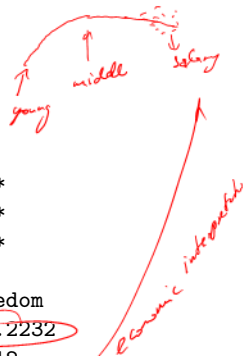
young · middle · salary

```
Coefficients:
               Estimate Std. Error t value Pr(>|t|)
(Intercept) 10.0419773  0.4559986   22.022  < 2e-16 ***
age          0.1731310  0.0238317    7.265 7.96e-12 ***
I(age^2)    -0.0019771  0.0002898   -6.822 1.02e-10 ***

Residual standard error: 0.5608 on 202 degrees of freedom
Multiple R-squared: 0.2308,     Adjusted R-squared: 0.2232
F-statistic:  30.3 on 2 and 202 DF,  p-value: 3.103e-12
```

economic interpretation

$\beta_2 < 0$          $\beta_2 > 0$

Universität Augsburg

*Cross-validation for bandwidth*

```
model.np = npreg(logwage ~ age, regtype = "lc", bwmethod = "cv.aic",
                 data = cps71)
summary(model.np)



Regression Data: 205 training points, in 1 variable(s)
                  age
Bandwidth(s): 1.551218

Kernel Regression Estimator: Local-Constant
Bandwidth Type: Fixed
Residual standard error: 0.2750934
R-squared: 0.3261299
```
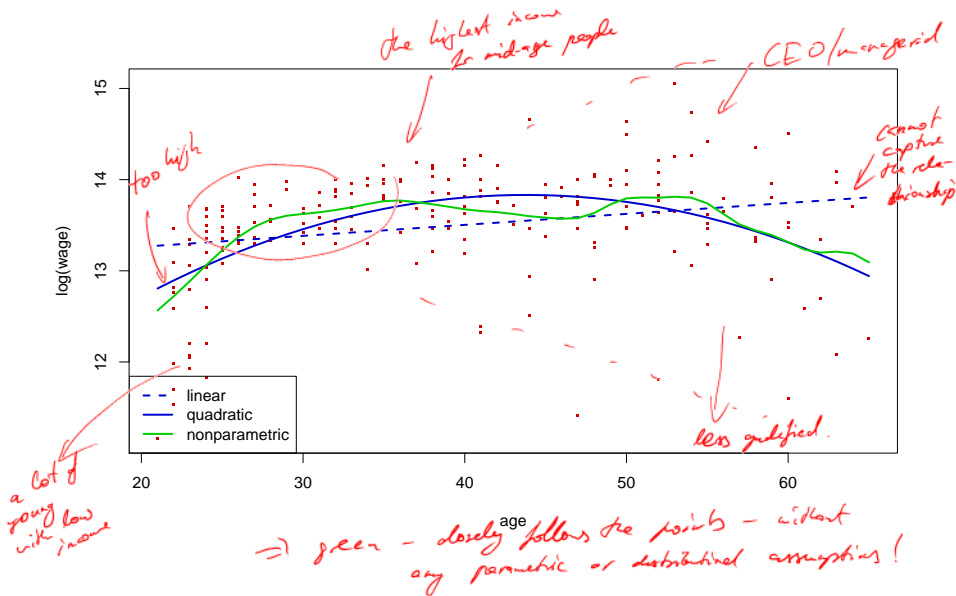
$\Rightarrow$ 10% than for parabola
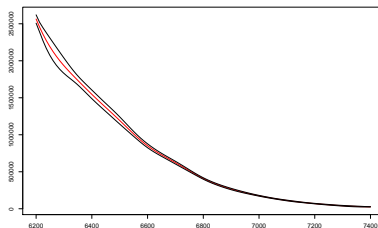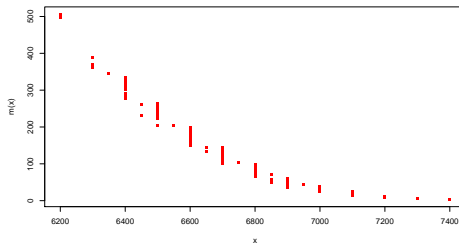
### Example: options

Consider call options on the DAX-certificates. The price of the option $C = C_t(S, K, \tau, r, \sigma^2)$ at time point $t$ depends in a complex way on the asset price $S$, strike price $K$, time to maturity $\tau$, risk-free rate $r$ and the volatility $\sigma^2$.

### Data:

Call options prices on 17.01.2001 with 1 month to maturity and with different strike prices.

$$C_i = m(K_i) + \varepsilon_i.$$

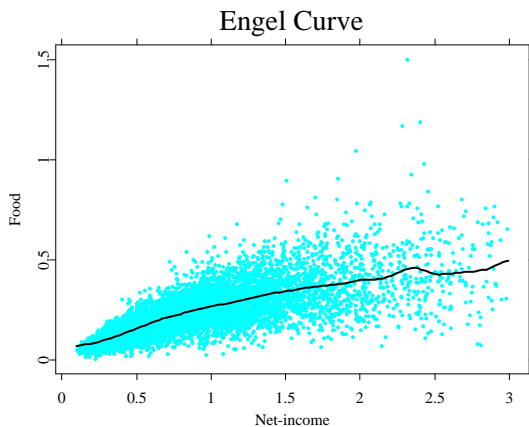# The NW estimation implies the following behaviour.

Figure: Nadaraya-Watson kernel regression, $h = 0.2$, U.K. Family Expenditure Survey 1973

KOE: $h_{opt}$ by minimizing MISE ⟹ Silverman's rule of thumb

## Bandwidth selection: Cross Validation

separate estimation and validation by using leave-one-out estimators

out-of-sample prediction for
the $i$-th observation

$$\min_{w.r.t. h} \quad CV(h) = \frac{1}{n} \sum_{i=1}^{n} \{Y_i - \widehat{m}_{h,-i}(X_i)\}^2 w(X_i) \quad \longrightarrow \text{weighted CV.}$$

estimated m function, if we drop
the $i$-th observation

minimizing gives $\widehat{h}_{CV}$

Out-of-sample forecasts

Important: not really correct for time series data!