# **Data 100**
# Lecture 7: EDA &
# *Visualization*

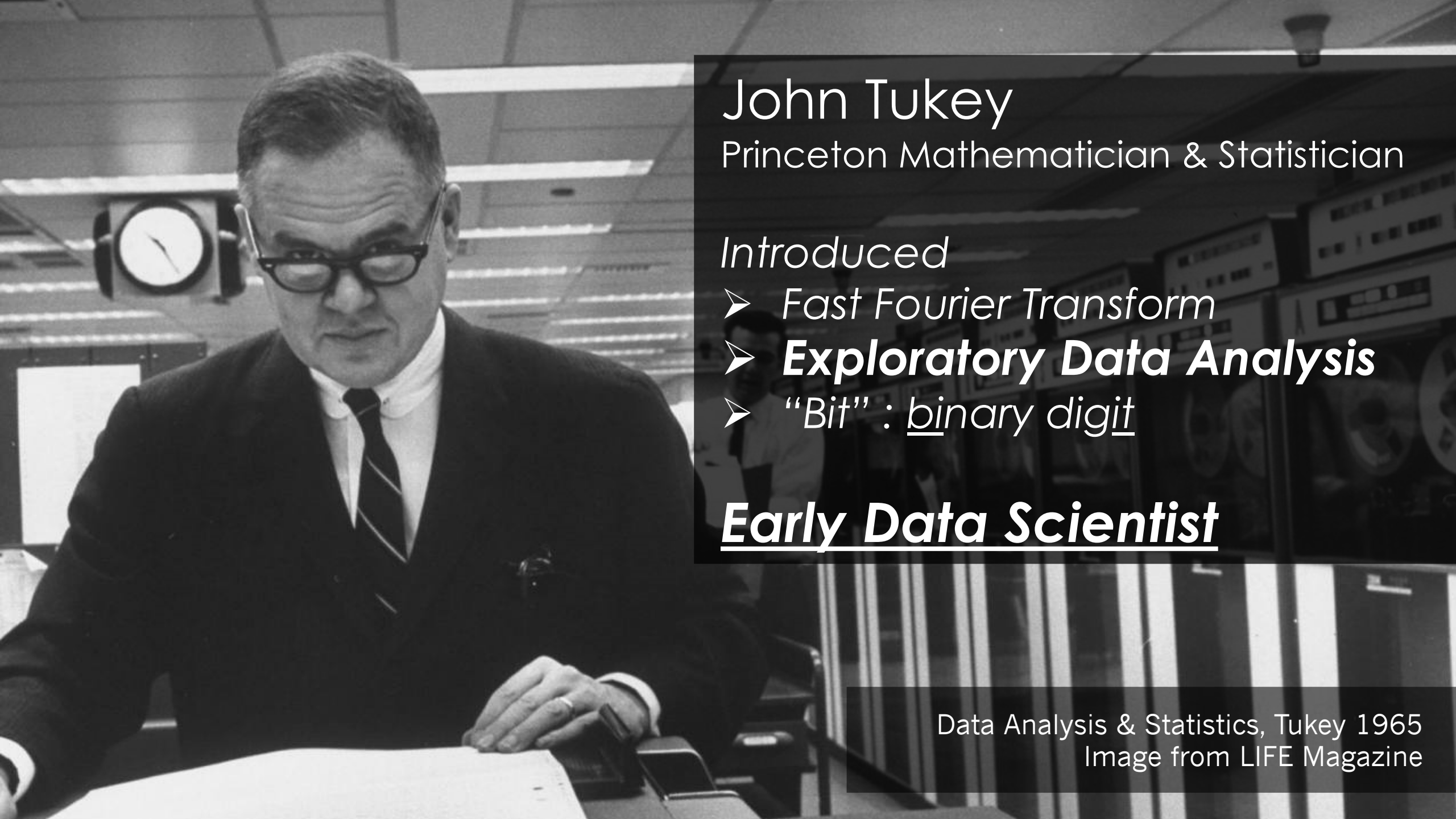# Exploratory Data Analysis (EDA)

*"Getting to know the data"*

A process of transforming, visualizing, and summarizing data to:

- ➢ Build/confirm understanding of the data
- ➢ Identify and address potential issues in data
- ➢ Inform the subsequent analysis
- ➢ Discover *potential* relationships

- ➢ **EDA is an open-ended analysis**
  - ➢ Be willing to find something surprising

# Exploratory Data Analysis (EDA)

*"Getting to know the data"*

➢ We used EDA with the $CO_2$ data and DAWN data to check the quality of the data.

➢ We also use EDA to help prepare for formal modeling.

➢ We also use EDA to confirm our modeling was reasonable

➢ Plots can uncover features, distributions, and relationships that can't be detected from numerical summaries
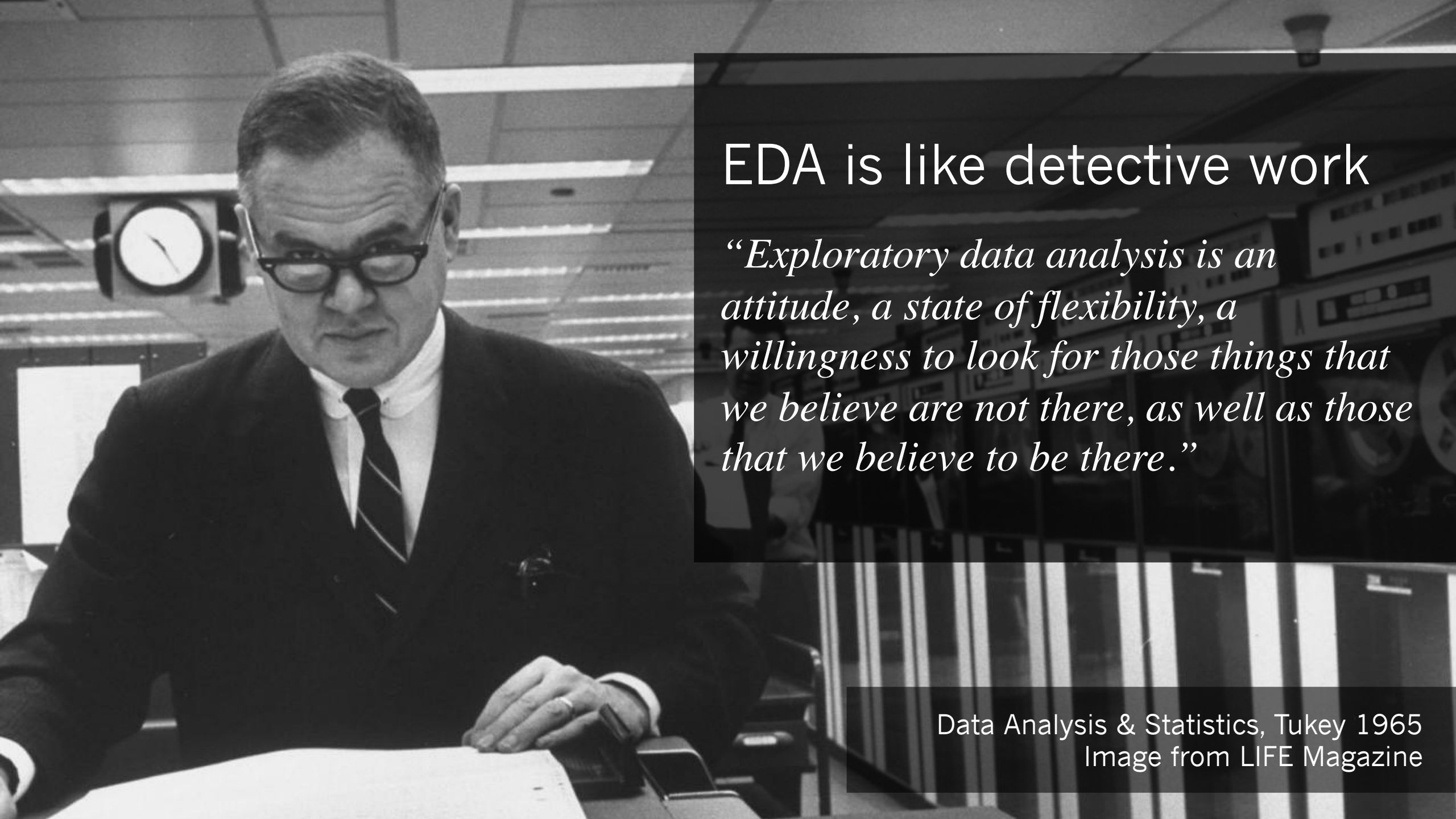
John Tukey
Princeton Mathematician & Statistician

Introduced
➤ *Fast Fourier Transform*
➤ **Exploratory Data Analysis**
➤ *"Bit" : binary digit*
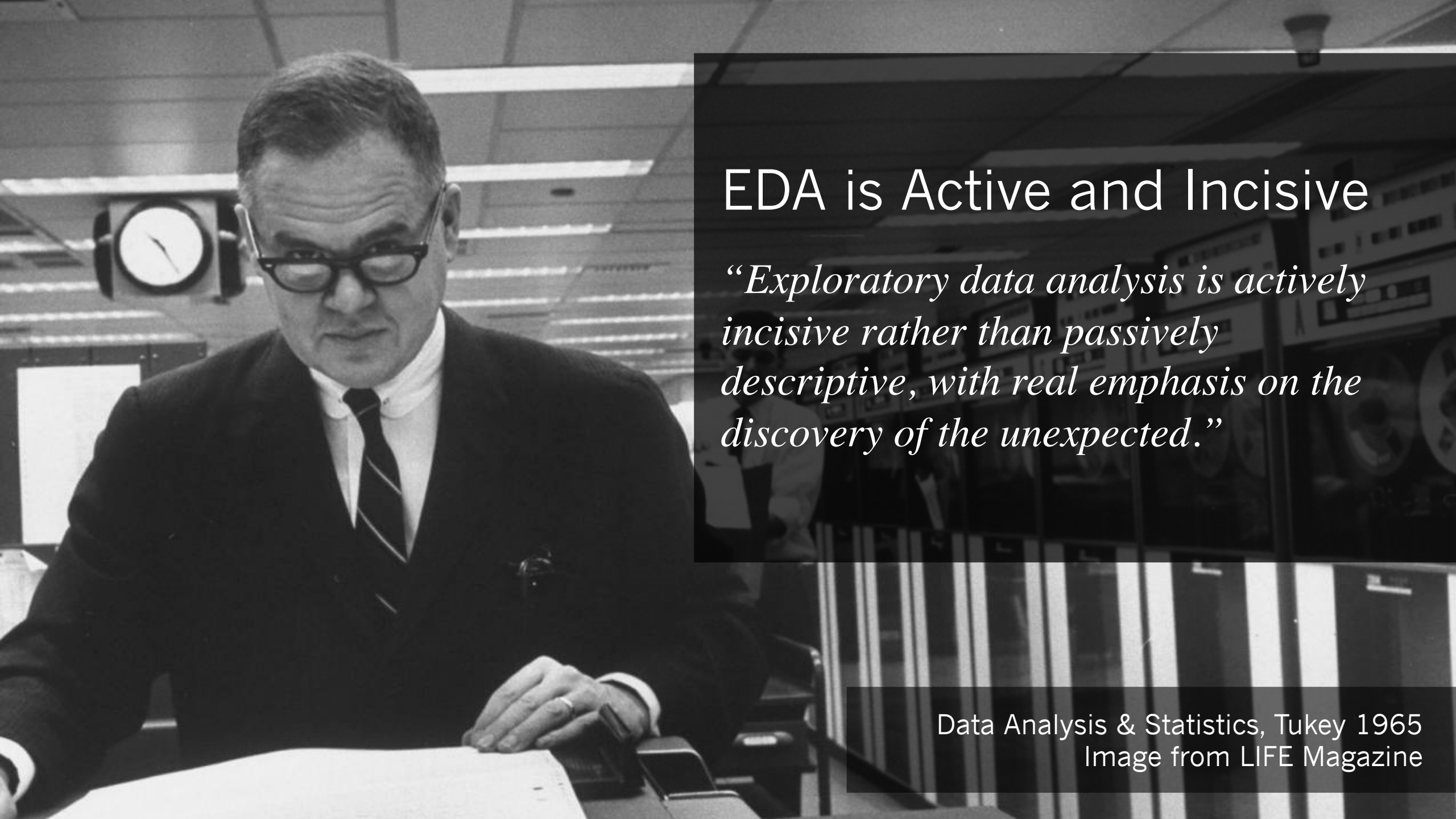
**Early Data Scientist**

Data Analysis & Statistics, Tukey 1965
Image from LIFE Magazine

# EDA is like detective work

*"Exploratory data analysis is an attitude, a state of flexibility, a willingness to look for those things that we believe are not there, as well as those that we believe to be there."*

Data Analysis & Statistics, Tukey 1965
Image from LIFE Magazine

# EDA is Active and Incisive

*"Exploratory data analysis is actively incisive rather than passively descriptive, with real emphasis on the discovery of the unexpected."*

Data Analysis & Statistics, Tukey 1965
Image from LIFE Magazine

# The Variable Represents

# The Variable Represents

Urban Dictionary:

Go and be a good example to the others of your group or in your position

Huh?

*A Variable represents a feature*

It is distinct from it's coding in a data file or data frame. It is more than a column in a table.

# Variable

## Quantitative

Ratios and intervals have meaning.

### Continuous

Could be measured to arbitrary precision.

**Examples:**
- Price
- Temperature

### Discrete

Finite possible values

**Examples:**
- Number of siblings
- Yrs of education

## Qualitative

### Ordinal

Categories w/ levels but no consistent meaning to difference

**Examples:**
- Preferences
- Level of education

### Nominal

Categories w/ no specific ordering.

**Examples:**
- Political Affiliation
- CalD number

| | Quantitative Continuous | Quantitative Discrete | Qualitative Nominal | Qualitive Ordinal |
|---|---|---|---|---|
| $CO_2$ level | X | | | |
| Number of siblings | | X | | |
| GPA | X | | | |
| Income bracket | | | | X |
| Race | | | X | |
| Number of years of education | | X | | |
| Yelp Rating | | | | X |
| Lane of traffic (left, middle, right) | | | X | |
| Left GRADE in 100 | | | | X |

# Basic Plots

Match Variable Type to Plot Type

# Basic Visualizations

➢ How to choose the "right" one(s)

➢ How to read them –
  ➢ Distributions
  ➢ Relationships

# Kaiser Study

➢ Oakland Kaiser mothers

➢ 1960s

➢ Measure the babies weight (in ounces) at birth

➢ All babies:
  ➢ Male
  ➢ Single births (no twins, etc.)
  ➢ Survived 28 days

# Information collected on mother's and their babies

➢ Birth weight (ounces) *Quantitative Continuous*

➢ Gestation (weeks)

➢ Parity - total number of previous pregnancies *Quantitative Discrete*

➢ Mother's height and weight

➢ Mother's smoking status

➢ Mother's age, race, education level, income level

➢ Father's information and more...

*Qualitative*

*Qualitative Ordinal*

# One Variable

What is the Distribution of the values of the variable?

# Quantitative – Continuous

➢ Birthweight

➢ The most basic visual representation of one quantitative variable is the *rug plot*

Hard to see much of the distribution with this rug plot



one thread for each observation

# Birthweight



Normalized birth weight distribution of babies

## Histogram

With the histogram we hide the details of individual observations and view the general features of the distribution.

How would we describe the distribution of birth weight?

# Distribution Features

*Normal Tails*

- Modes
    - Number  1
    - Location  near 120 oz
    - Size  main mode

- Symmetry
    - Symmetric  very symmetric
    - Skewed left or right  slight left skew

- Tails
    - Long, short, "normal"

- Gaps

- Outliers

Normalized birth weight distribution of babies

Birth weight in ounces

# Distributions & Smoothing

# A Small Dataset

10 values
0.7, 0.8, 0.9, 2.1, 2.2, 2.8, 2.9, 3.1, 3.6, 4.8
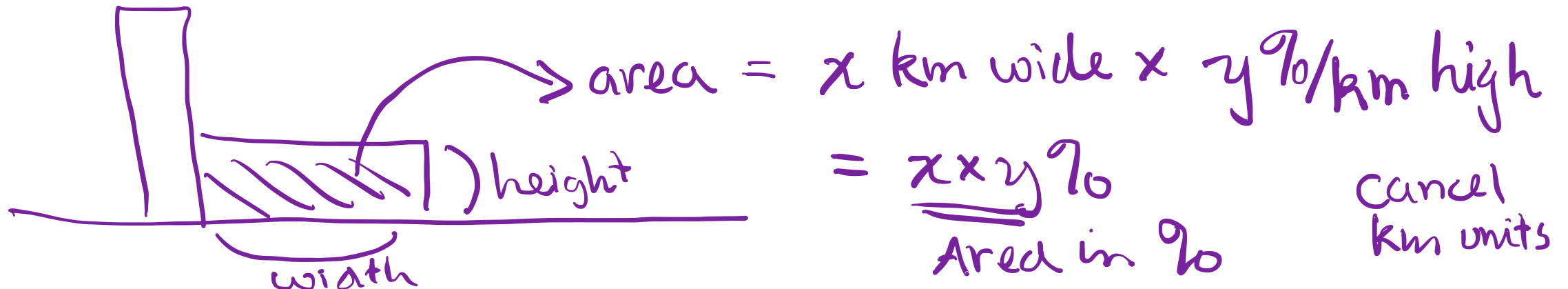
Rug Plot

Shows the location of each value

# We want to smooth these rug threads



BECAUSE

➢ this is a sample and we believe that other values near the ones we observed are reasonable

➢ we want to focus on general structure rather than individual observations

# Important Properties of Histograms

➢ Total Area of the bars = 100% (or 1)

➢ Units on the y-axis are percent/x-unit

➢ Area of a bar = percentage of values in that bar

unit matching:   x-units km      y-units %/km

area = x km wide × y %/km high

= x × y %

Area in %

cancel km units

width

height

# Example – One large bin from 0 to 5

one 'mode'
no/short tails

The 10 points are spread evenly across one large bin

fraction per km

$Area = 5 km \times 0.2 \ frac/km$

$= 1 \ fraction$
all observations

Not Very Informative Distribution
km

# Example



Bin width $1/4$ km

With these narrow bins, the histogram is little more than a rug plot

Area $= \frac{1}{4}$ km $\times$ 0.4 Frac/km

$= 0.1$ fraction of sample

Area $= \frac{1}{4} \times 0.8 = 0.2$ (2 observations)

# Example



Bins can be different widths

$$\text{Area} = 2 \text{ km} \times 0.15/\text{km}$$
$$= 0.3$$

km

# A histogram smooths

We want to smooth out these points because:



- ➢ this is a sample and we believe that other values near the ones we observed are reasonable

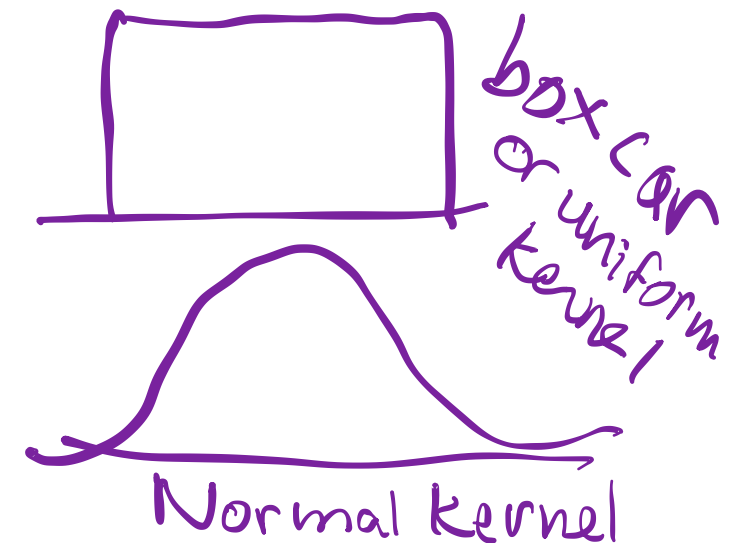- ➢ we want to focus on general structure rather than individual observations

The values 3.1, 3.6, and 4.8 have their proportion (3/10) spread over the bin [3,5] That is, without the rug, we can't tell where the points are in the bin
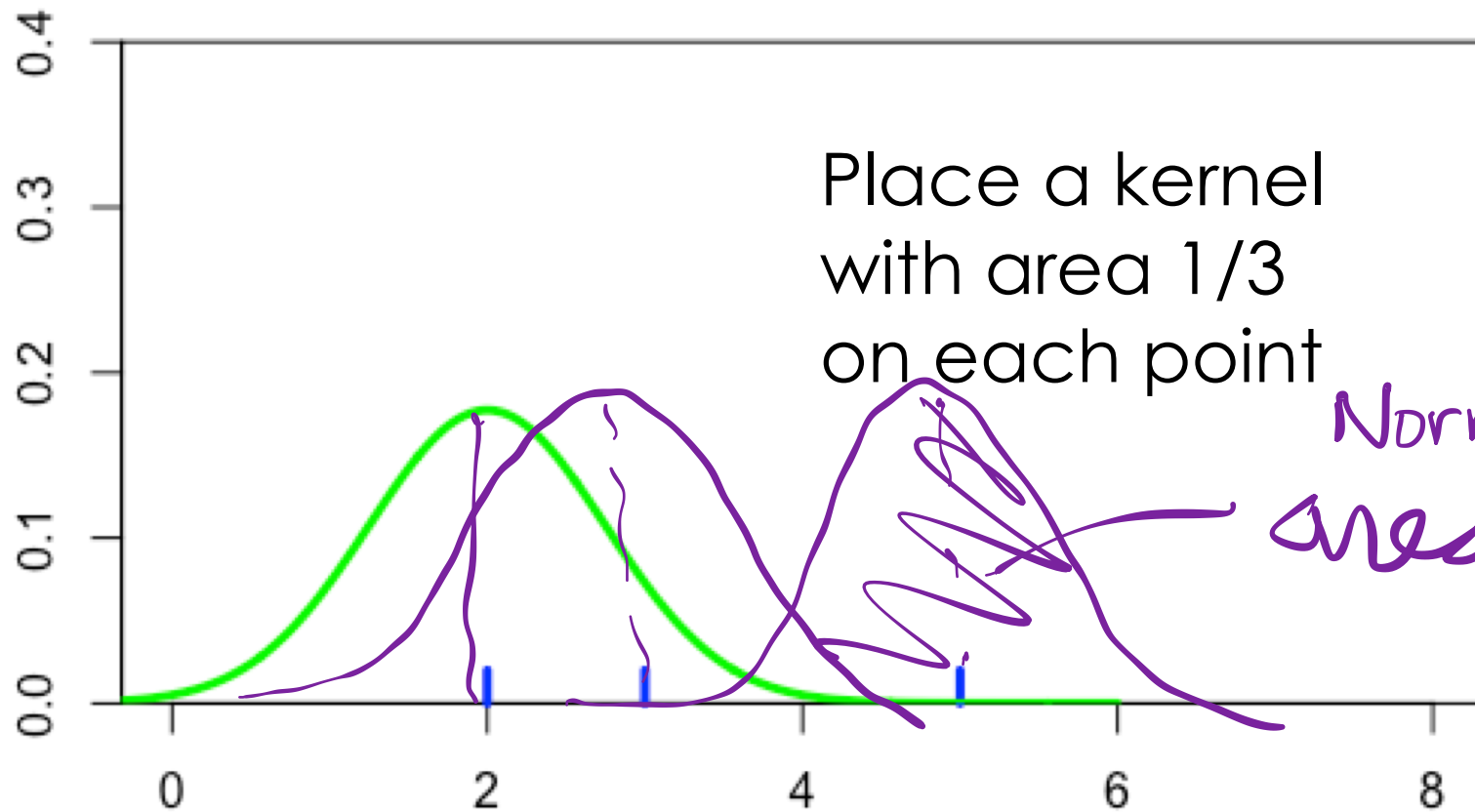
# Kernel Density Estimate: Alternative Smoother

Consider one point

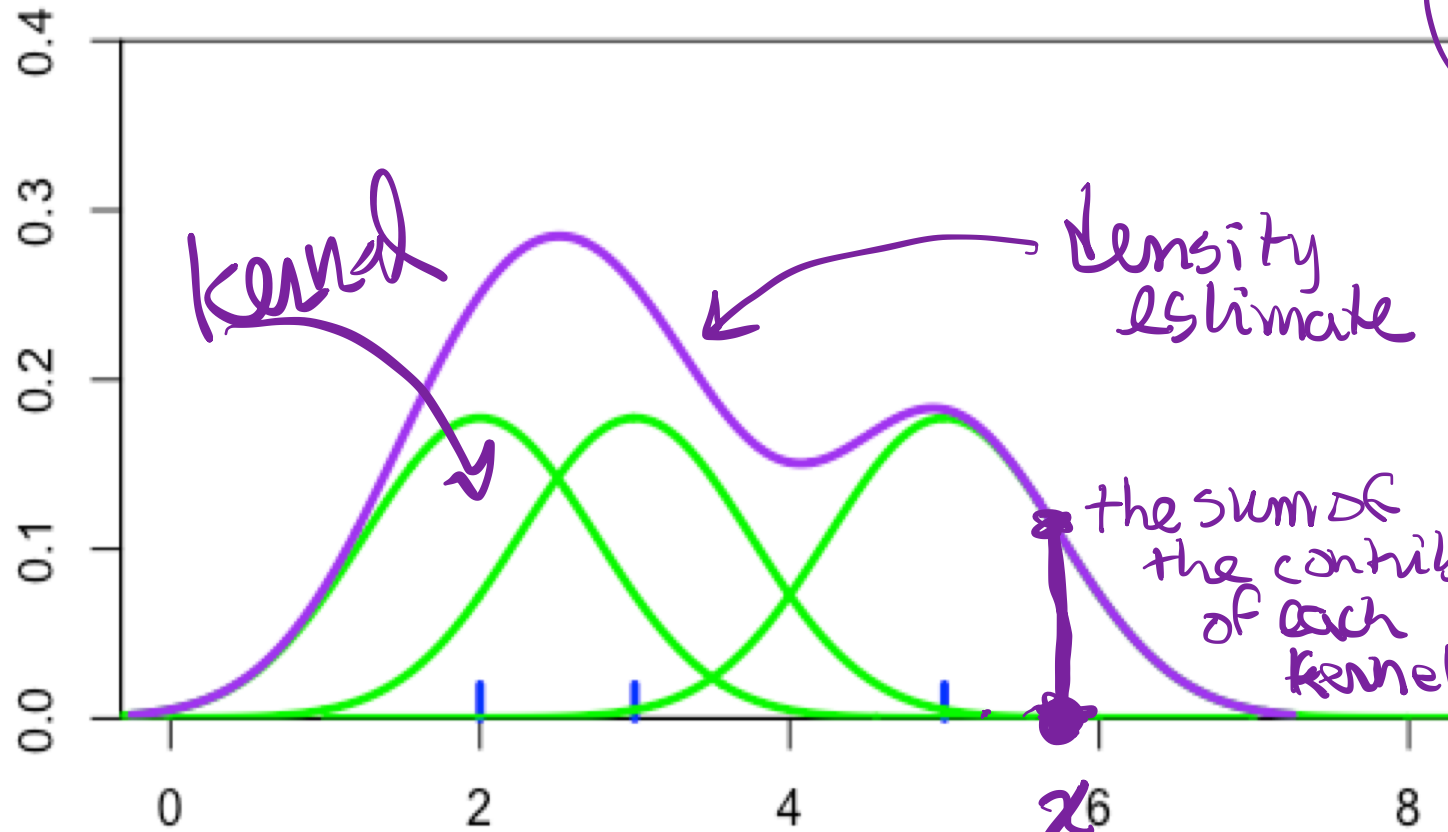Smooth with a kernel function, rather than in a histogram bin



Area under this curve is 1 or 100%

kernel function

∠ Area

box car or uniform kernel

Normal kernel

# 3 points –
# each represents 1/3 of the data



Place a kernel
with area 1/3
on each point

Normalize the kernels by the number of observations

area 1/3

# KDE – 3 points
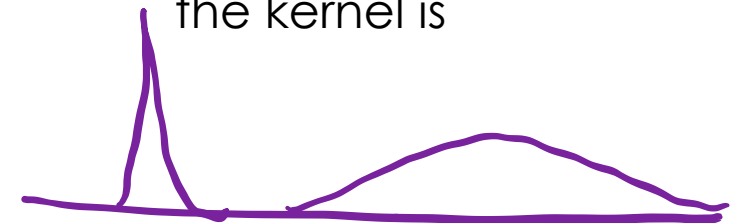
Sum the 3 kernels at each point to get the density curve

#obs

$$f(x) = \frac{1}{n} \sum_{i=1}^{n} K_h(x - x_i)$$



kernel

density estimate

the sum of the contribution of each kernel function

$x$

$K_h$ is the green Kernel function

h refers to how peaked /spread the kernel is

Example    Flat kernel

Density Curve
is akin to **1** bar
histogram

Broad Flat
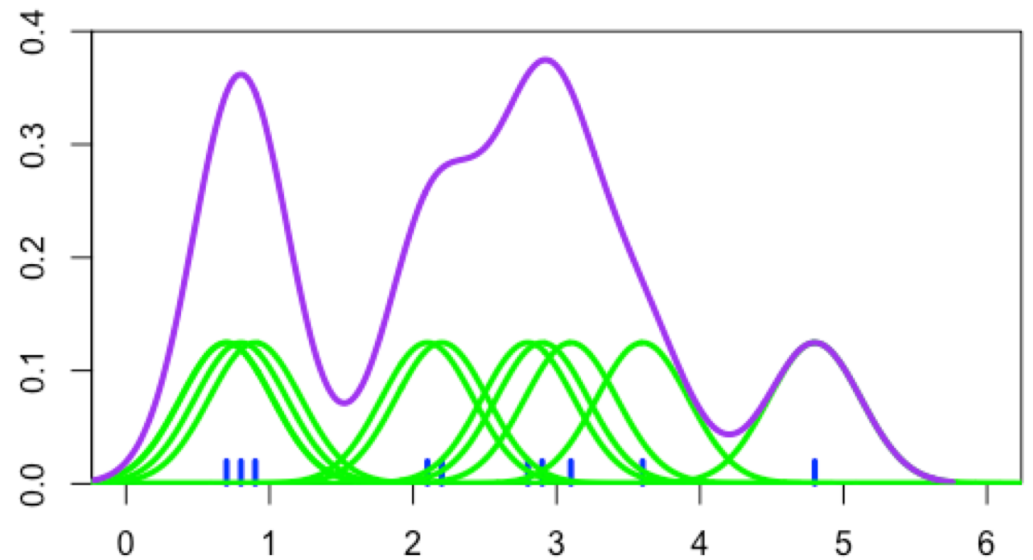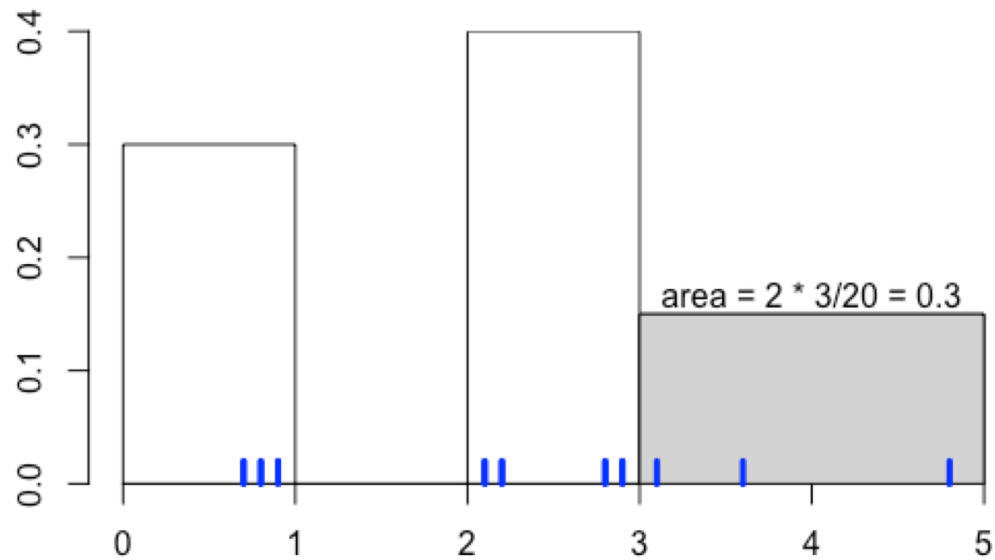Density Curve

kernel centered at

# Example



Density Curve is too grainy
It's like a rug plot

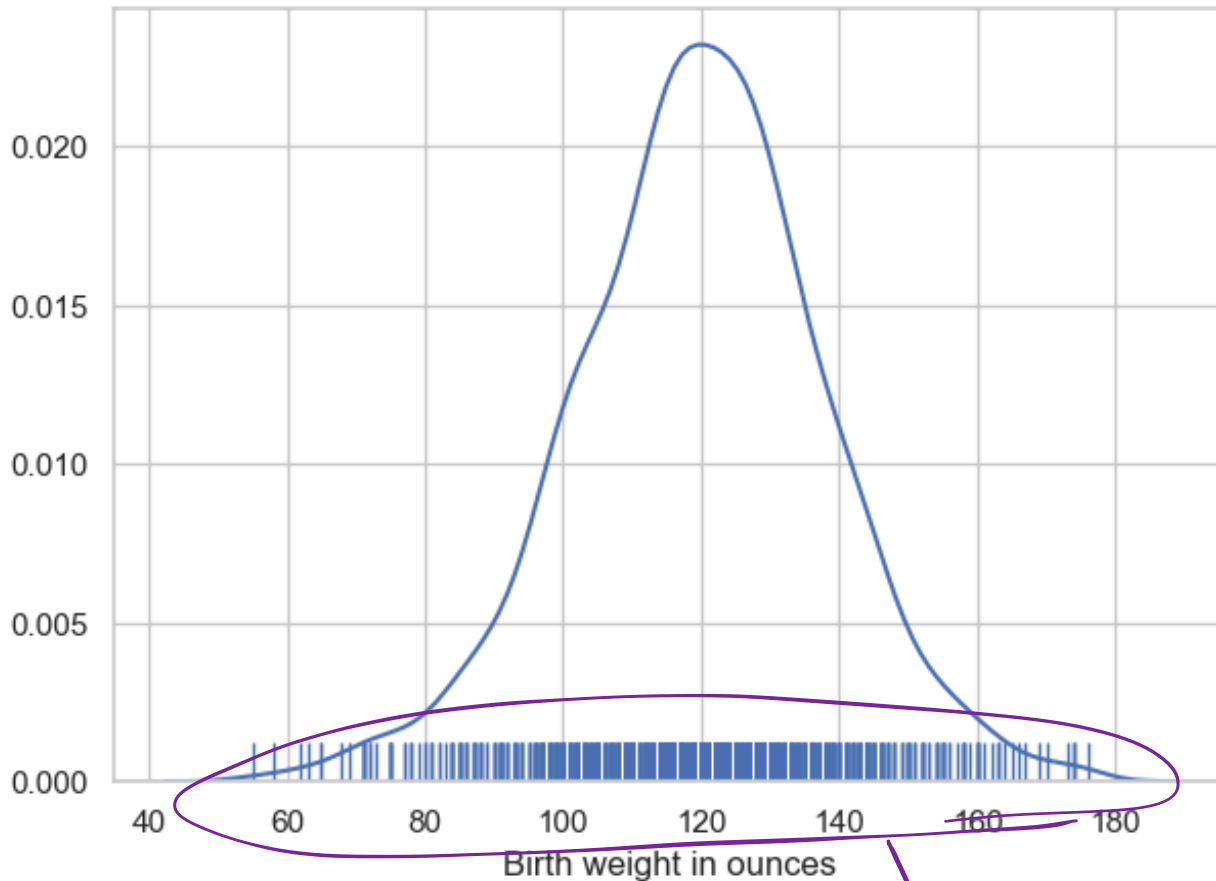The software chooses a kernel bandwidth for you, but you can also specify your own.

# Compare the Histogram and the KDE

0.7, 0.8, 0.9, 2.1, 2.2, 2.8, 2.9, 3.1, 3.6, 4.8



area = 2 * 3/20 = 0.3

# Birthweight – Density Curve



Normalized birth weight distribution of babies

rug plot not very informative

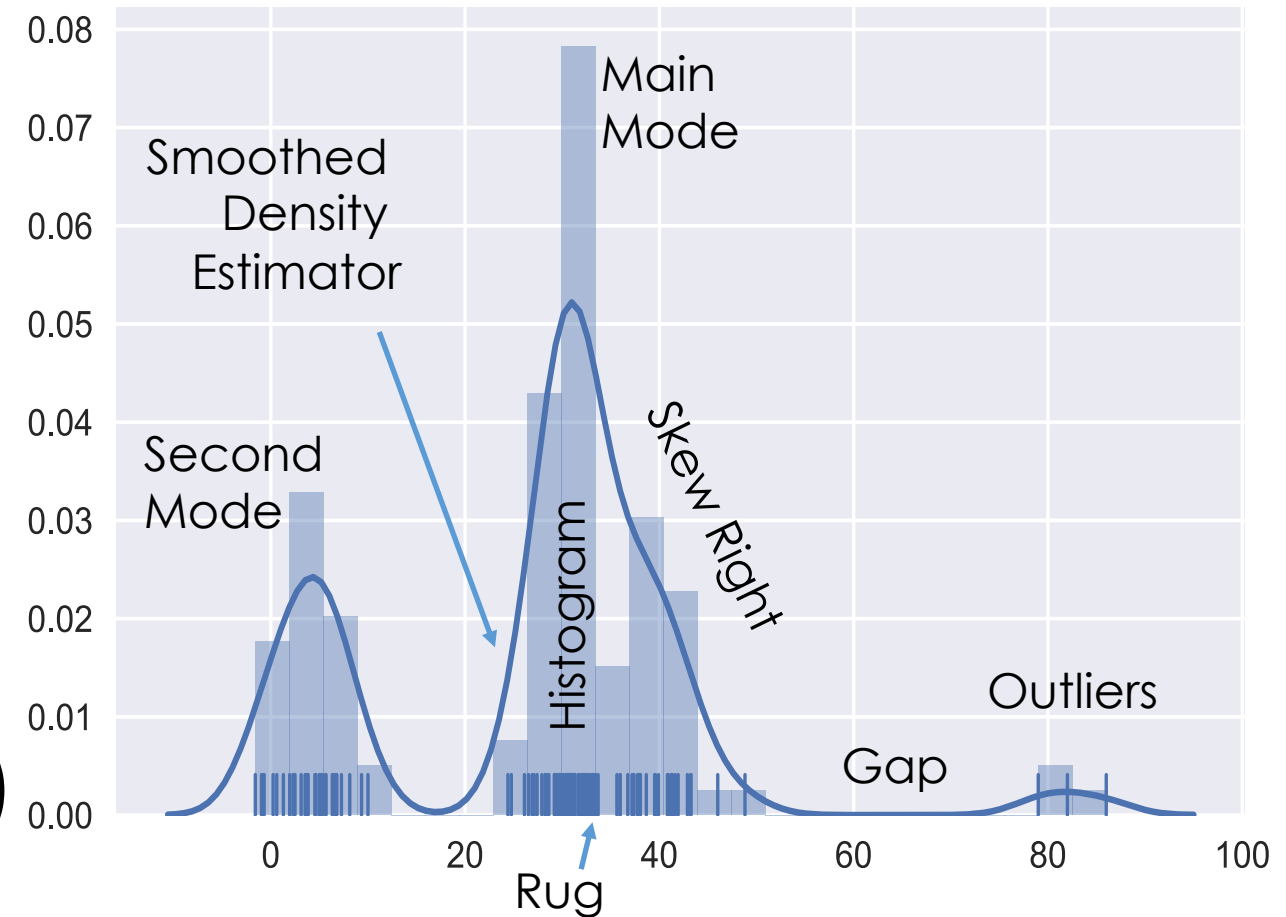How would we describe the distribution of birth weight?

Unimodal

Main mode at 120 oz

Slight left skew

Tails about normal

# Histograms and Density Curves

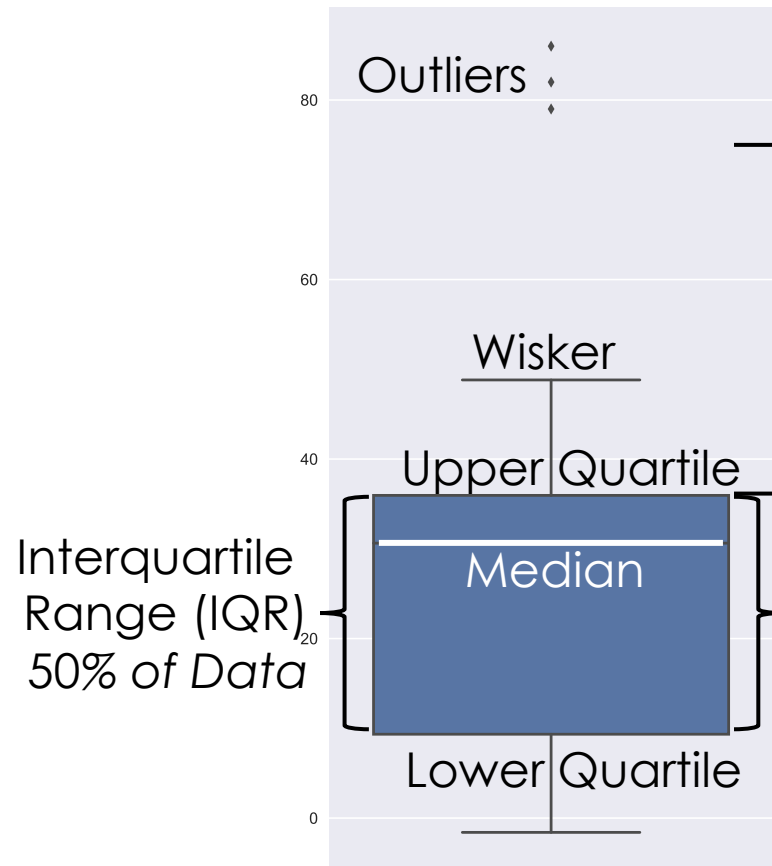Describes distribution of data – relative prevalence of values

➢ Histogram
  ➢ relative frequency of values
  ➢ Tradeoff of bin sizes

➢ Rug Plot
  ➢ Shows the actual data locations

➢ Smoothed density estimator
  ➢ Tradeoff of "bandwidth" parameter (more on this later)

# Box Plot

➢ Useful for summarizing distributions and comparing multiple distributions

Outliers

80

Wisker

60

Upper Quartile

40

Median

Interquartile
Range (IQR)
50% of Data

20

Lower Quartile

0

Outliers are more than 1.5 * IQR away from lower and upper quartiles.

Visualization of summary statistics

Can lose a lot of features, such as…?

Modes
Gaps

# Our Data

| 1 | 2 | 3 | 4 | 5 | 6 | 4 | 3 | 2 | 1 |

| | | | | 6 | 5 | | | | |

0.7, 0.8, 0.9, 2.1, $\boxed{2.2, 2.8}$ 2.9, 3.1, 3.6, 4.8

- ➢ Median
- ➢ Lower Quartile
- ➢ IQR
- ➢ Hinge

Tukey's short cut for finding these values

$n = \#obs = 10$

median is the $\frac{n+1}{2}$ smallest or largest obs

$= \frac{10+1}{2} = 5.5$ average the $5^{th}$ & $6^{th}$

median $= \frac{2.2 + 2.8}{2} = 2.5$

# Our Data

$$\overset{1}{0.7}, \overset{2}{0.8}, \overset{3}{\boxed{0.9}} \; 2.1, 2.2, 2.8, 2.9, \overset{3}{\boxed{3.1}} \overset{2}{3.6}, \overset{1}{4.8}$$

➢ Median

➢ Lower Quartile

➢ IQR

➢ Hinge

To find the quartiles, we take the count-down value for the median, drop the ½ (if it has it), and add one & divide by 2, e.g.)

$$5.5 \longrightarrow \frac{5+1}{2} = 3$$

IQR is UQ−LQ
$$= 3.1 − 0.9 = 2.2$$

Hinge is 1.5 IQR = 1.5 × 2.2 = 3.3

LQ is 3 in from bottom

UQ is 3 in from top

Any value more than 3.3 away from LQ/UQ is an outlier

# Quartiles from Tukey's "depth"

➤ Depth of the Median = (n + 1)/2
  ➤ Count in from top or bottom of ordered set of values
  ➤ If depth has a half then average the two values on either side

➤ Depth of Quartile = (round(m) + 1)/2
  ➤ Round the median depth down to nearest integer
  ➤ Count in from bottom to get the LQ and from the top to get the UQ
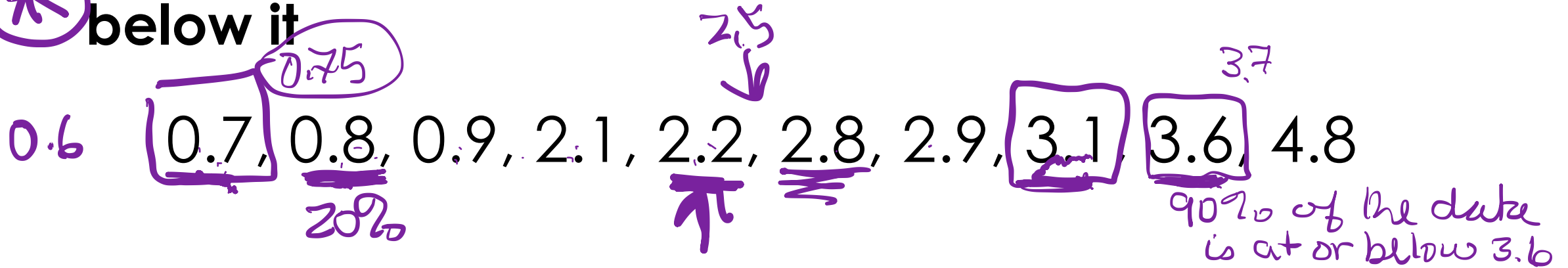  ➤ If depth has a half in it then average the two values on either side

# Percentile – Need a more general def

Notice the percentile will always correspond to a data point

➤ The P[th] percentile of a set of data is:

✳ **Smallest** value that has **at least** P% of the data **at or below it**

2.5

0.75

3.7

0.6   0.7, 0.8, 0.9, 2.1, 2.2, 2.8, 2.9, 3.1, 3.6, 4.8

20%

90% of the data is at or below 3.6

10th%tile = 0.7        90th%tile = 3.6        60th%tile = 2.8

15th%tile = 0.8        83rd%tile = 3.6        66th%tile = 2.9

Any value below 3.6 will not have 83% below it

# Percentile – with weighted data

➢ The P[th] percentile of a set of data is:

**Smallest** value that has **at least** P% of the data **at or below it**



sorted wages

5. 5. 5. 5. 5. 5. 20. 20. 20. 20. 50. 50.

corresponding weights

½  ½  ½  ½  ½  ½  1  1  1  1  ½  ½

sum to 8

$\frac{1}{2} + \frac{1}{2} + \frac{1}{2} + \frac{1}{2} + \frac{1}{2} + \frac{1}{2} + \frac{1}{2}$

50th %tile = **20**

4/8 = 50% is at or below 20

75th %tile = 20

7/8 = 87.5% is at or below 20

3/8 = 37.5 is at or below 5

> nothing in between

# Quantitative Discrete

We look for the same features

➢ Symmetry and skew

➢ Modes (number, location, and size)

➢ Tails (long, short, normal)

➢ Gaps

➢ Outliers

# Discrete Quantitative                # of Siblings



Bar plot – height not area is the proportion

What's the difference between these 2 plots?

# Qualitative

We look at the relative size of groups

➢ Equally distributed

➢ Symmetry, Modes, Tails and Gaps don't make sense

➢ Do most fall in one group?

Answers have implications in building prediction models

# Qualitative Variable

## Bar Width – has no meaning



Why do we not reorder the bars according from shortest to tallest?

# Education level

## Dot plot focuses on comparison of the values

*Per*

*Cleveland*

# Pairs of Variables

Combinations:

Both qualitative,

One qualitative and one Quantitative,

Both Qualitative

# Plotting Pairs of Quantitative Variables

➢ Scatter plot uncovers form of relationship between 2 variables

➢ Linear relationships are particularly simple to interpret

➢ Simple and elegant statistical theory for linear relationships

➢ Models are typically approximations, choose a simpler model over a complex one

# Common Relationships

**simple linear**

**simple nonlinear**

Ideal

Can also have more spread

Good too

Typically we transform to a linear rel.

**unequal spread**

**complex nonlinear**

Still linear but we need to take care when modeling

Very difficult to work with

height measured to nearest inch so we get these stripes

The scatter plot is a 2-d rug plot

mom's ht

# Hex Bin

Shading corresponds to density of points in the cell

histogram of each variable

180
160
140
bwt 120
100
80
60

55    60    65    70

Why hexagons —
Easier to see elliptical / linear relationships
More efficient for covering region
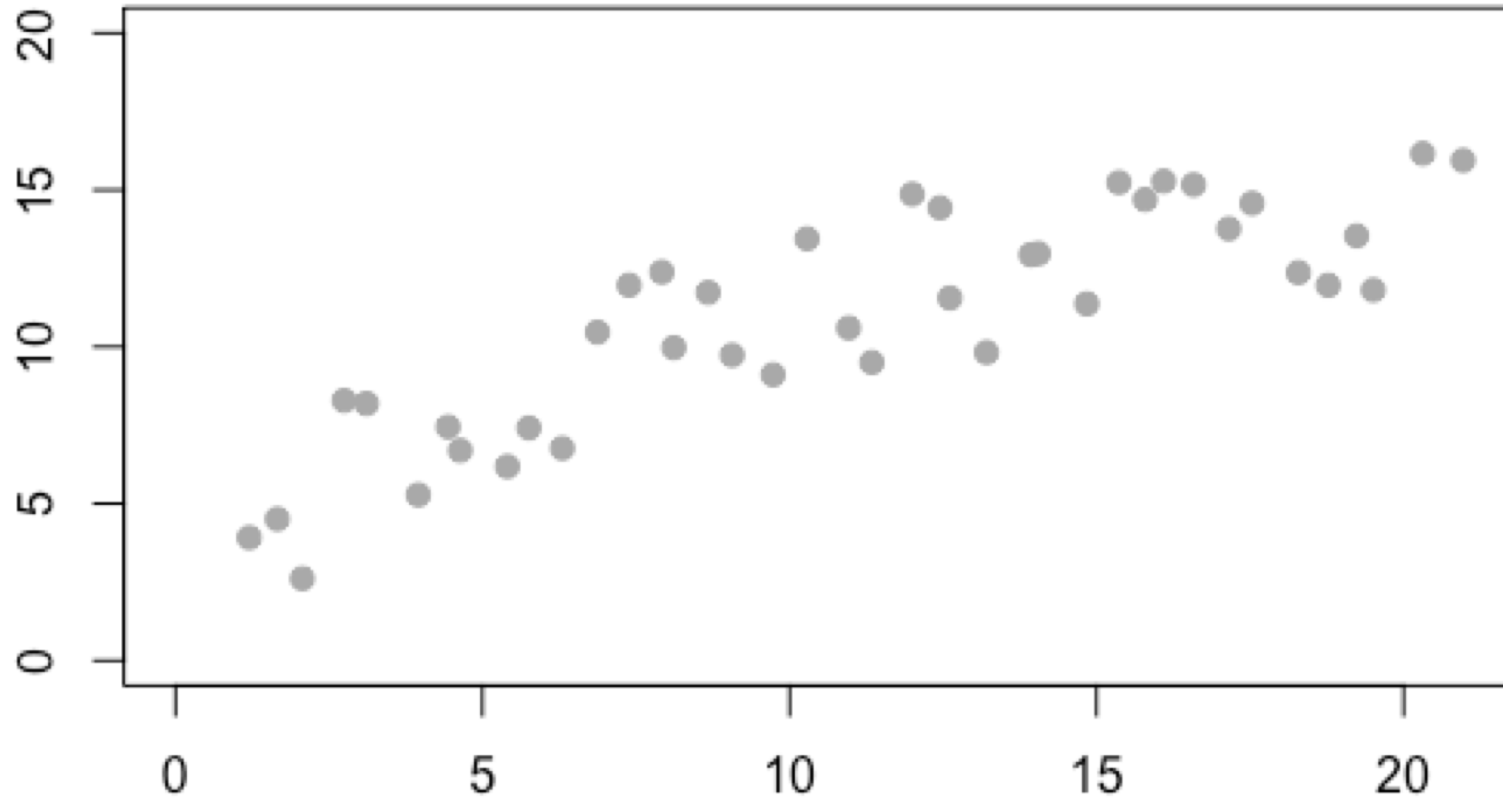Visual bias of squares — drawn to see vert & horizontal lines

# Smooth Contour

density curves for each variable

kde
in 2-d

Contours of
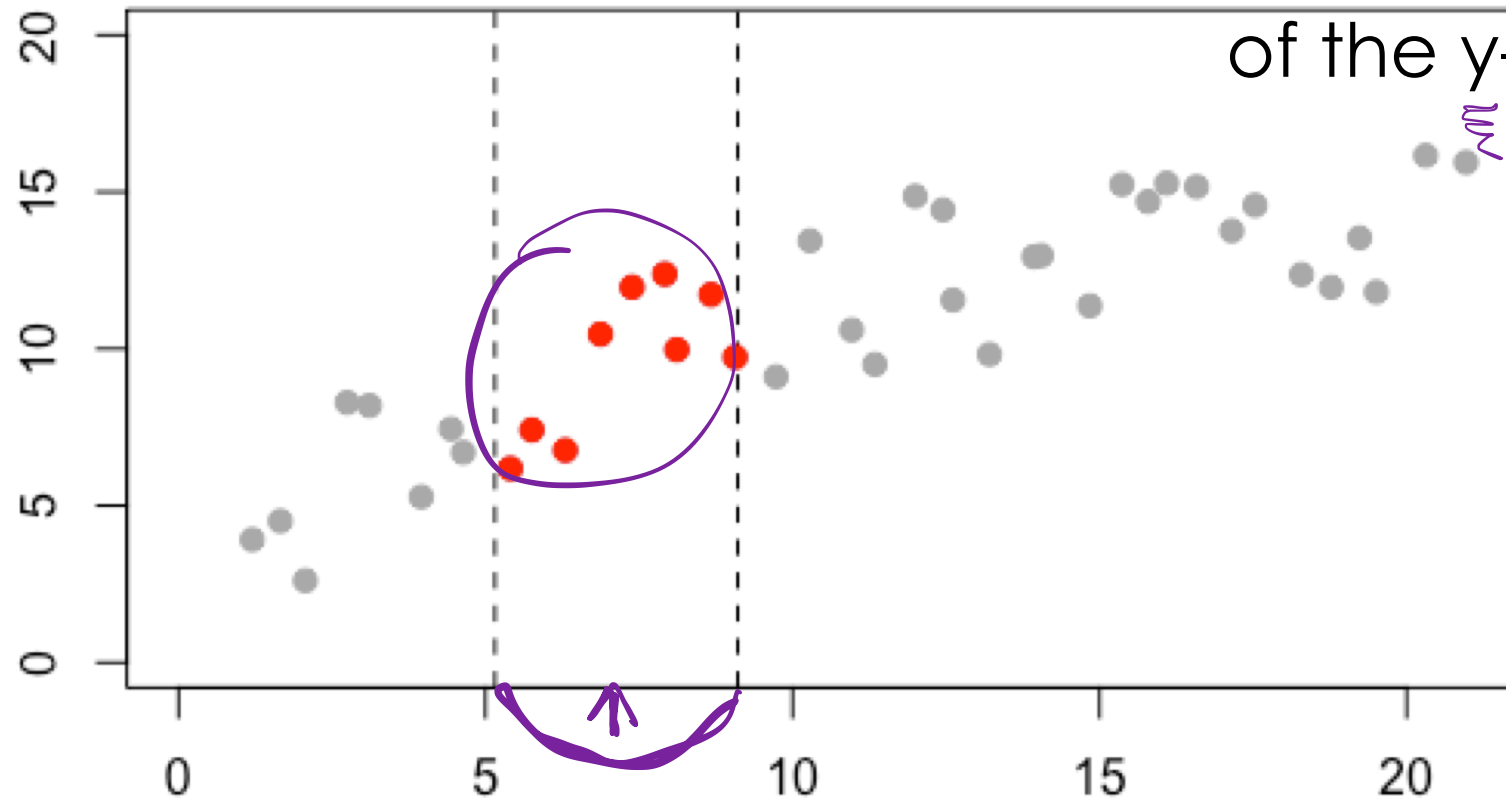the 3-d
density
smooth of
our 2-d data

# Smoothing Scatter plots

Now we want to smooth the y-values as a function of x
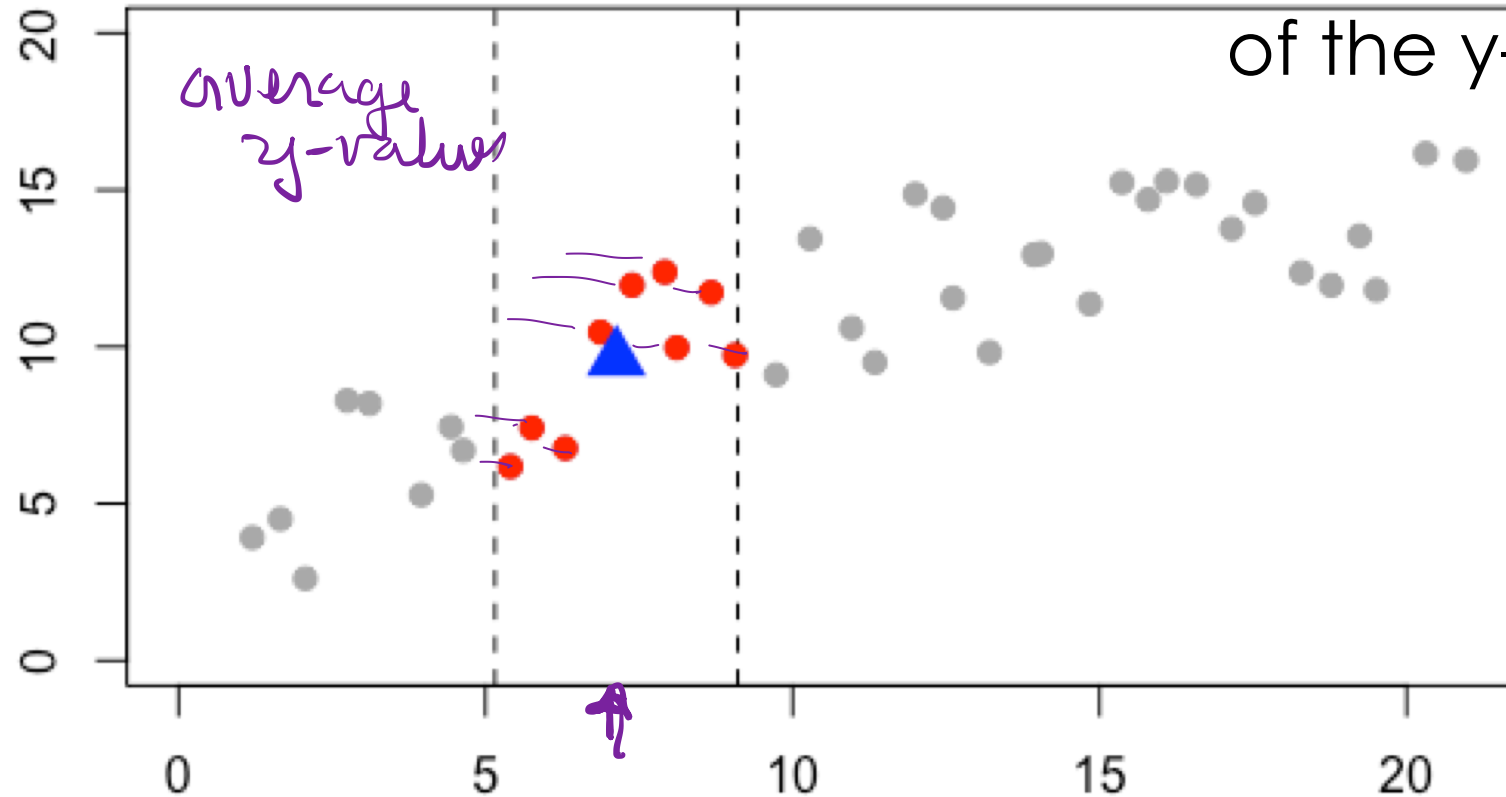
# Smoothing Scatter plots

For an x-value consider all of the x's near it
Take an average of the y-values
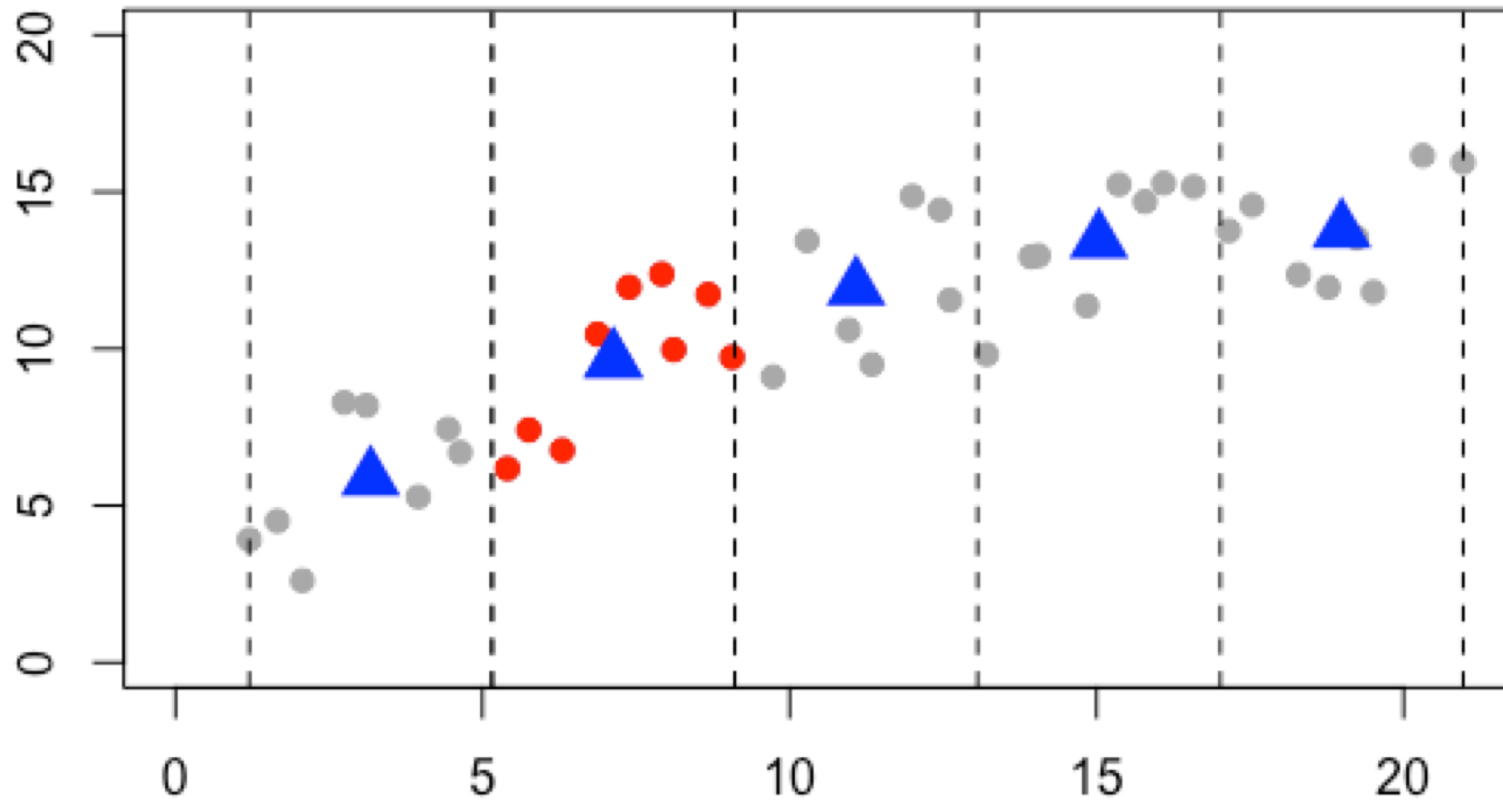


x-values
near our point of interest

# Smoothing Scatter plots

For an x-value consider all of the x's near it Take an average of the y-values
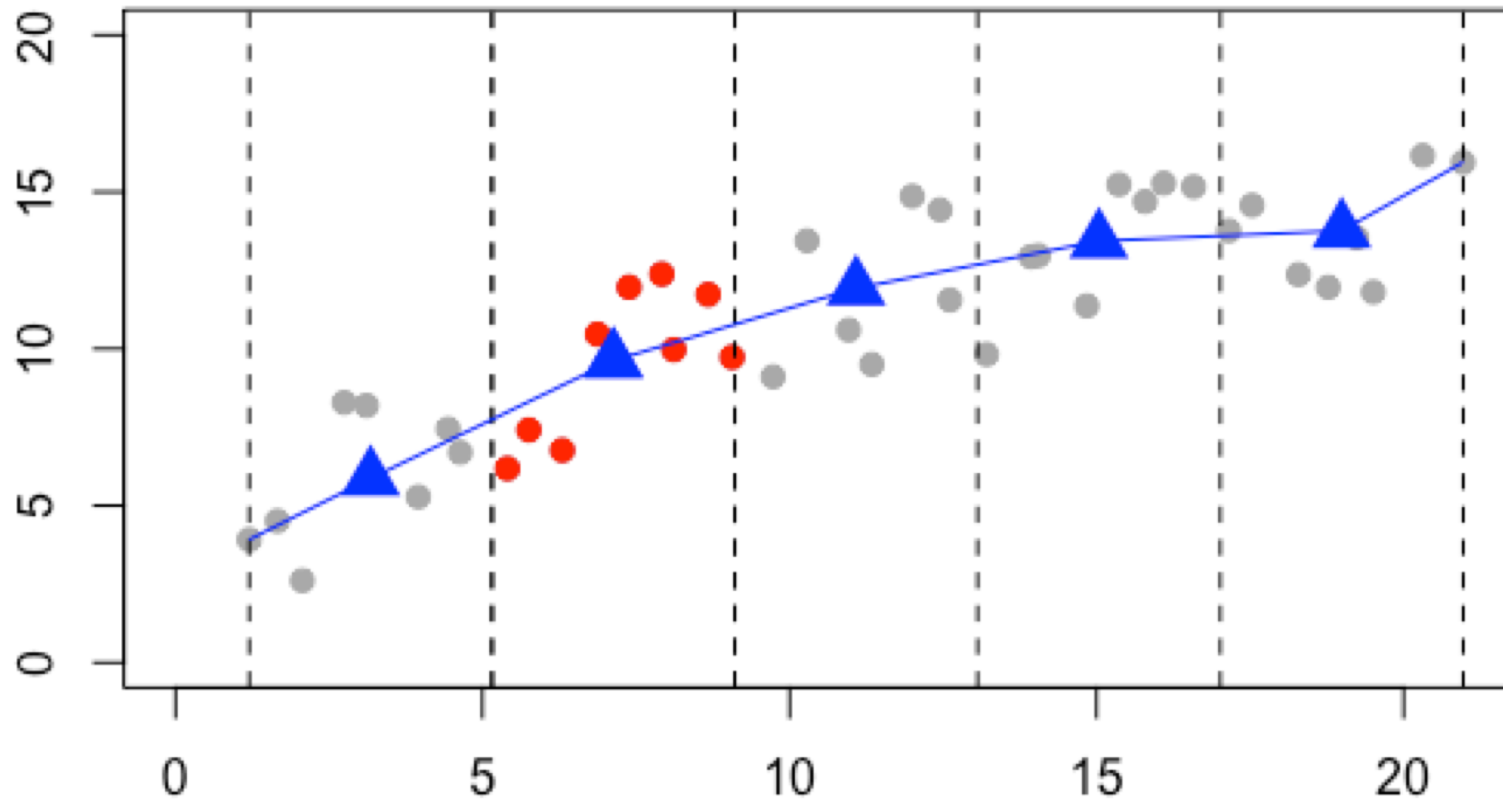
# Smoothing Scatter plots

Create bins for all x
Average y-values
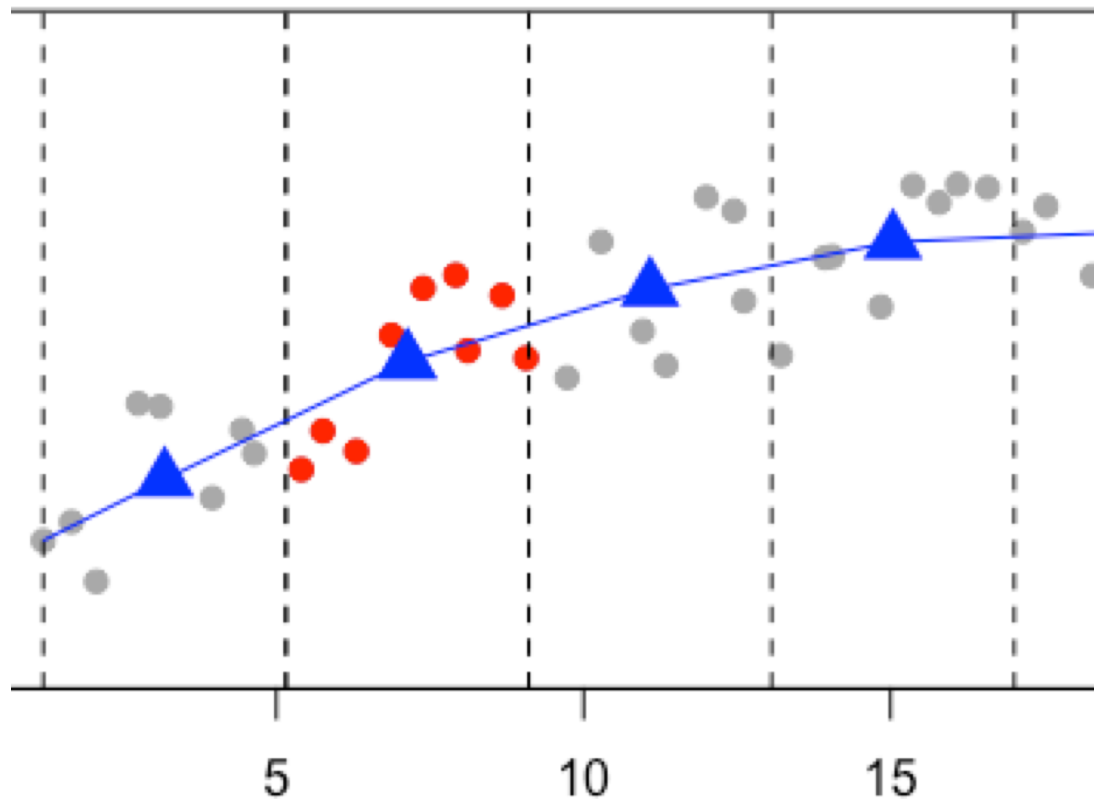in each bin

# Smoothing Scatter plots

These averages sketch out a curve



Connected these blue 🔺

# Smoothing Scatter plots



Rather than a simple average in fixed bins

We use kernels positioned on the $x_i$ to determine the weights to place on the $y_i$ in the average

_distance from $x$_

$$g(x) = \sum_{i=1}^{n} \frac{K_h(x - x_i)y_i}{\sum K_h(x - x_i)}$$

We average the y-values

The denominator ensures the weights sum to 1

Like with histograms we average based on distance from $x$

# Smoothing Scatter plots
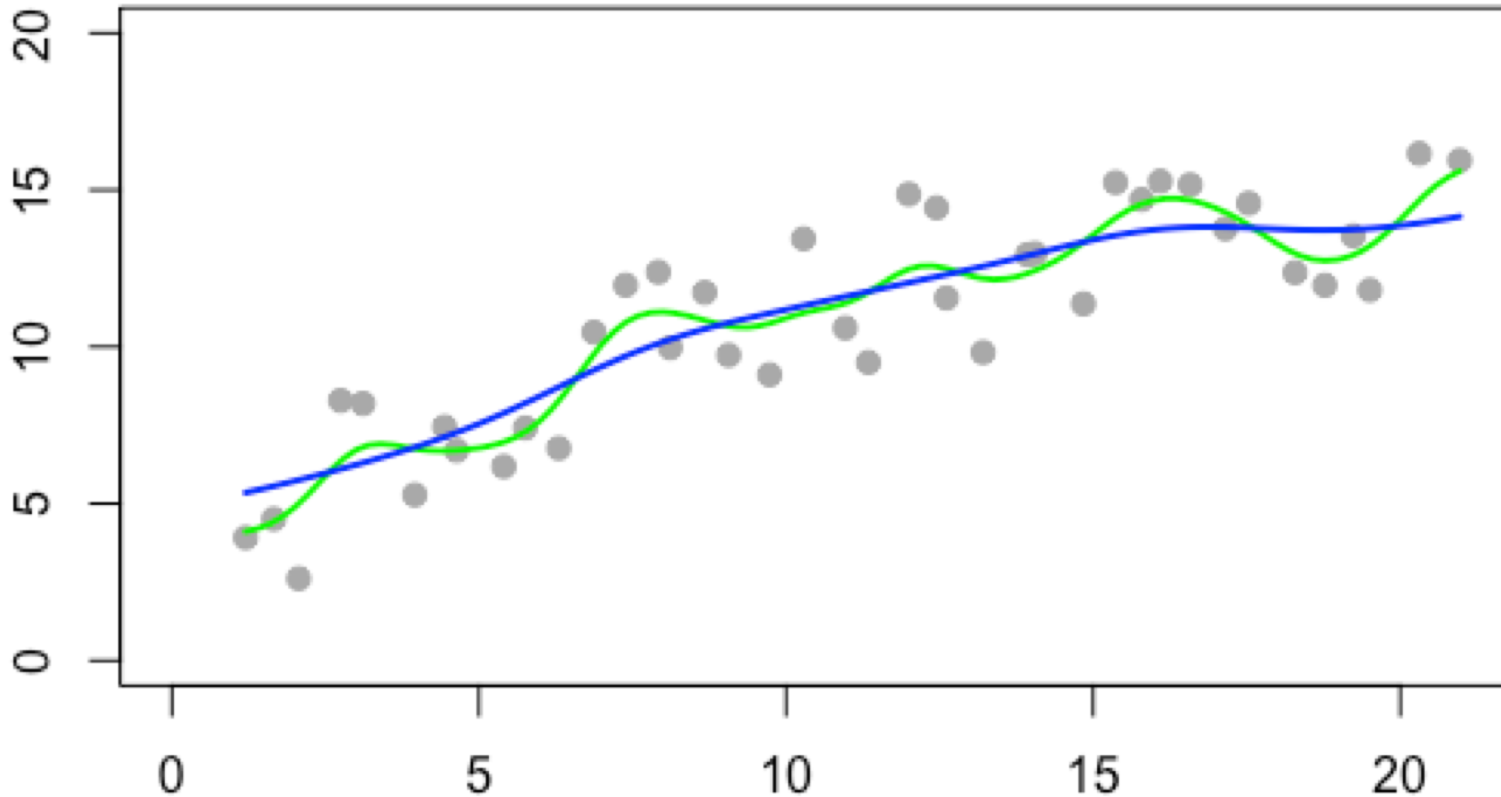
Rather than a simple average

We use kernels positioned on the $x_i$ to determine the weights to place on the $y_i$ in the average

$$g(x) = \sum_{i=1}^{n} \frac{K_h(x - x_i) y_i}{\sum K_h(x - x_i)}$$

The denominator ensures the weights sum to 1

# Smoothing Scatter plots

$$g(x) = \sum_{i=1}^{n} \frac{K_h(x - x_i) y_i}{\sum K_h(x - x_i)}$$



For each x, we find g(x) by a weighted average of the $y_i$

The $y_i$ are weighted according to the kernel function.
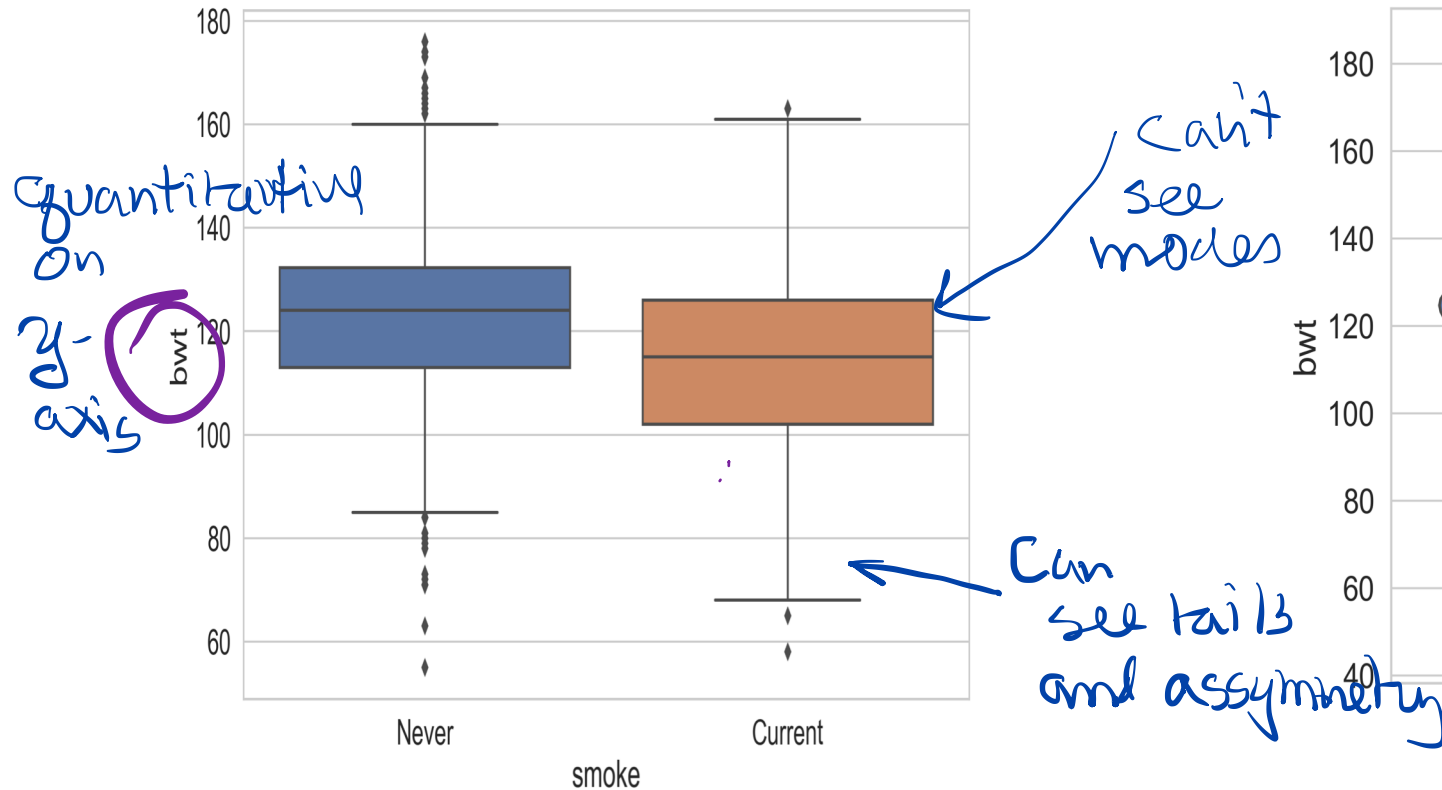So $x_i$ far from x do not contribute much to g(x)

# Local Smoothing

➤ Moving window

➤ Smooth/Average y values in the window

➤ Many different approaches for doing this:
  ➤ kernel methods (what we just showed),
  ➤ cubic splines, thin plate splines,
  ➤ Locally weighted smooth scatterplot (lowess)

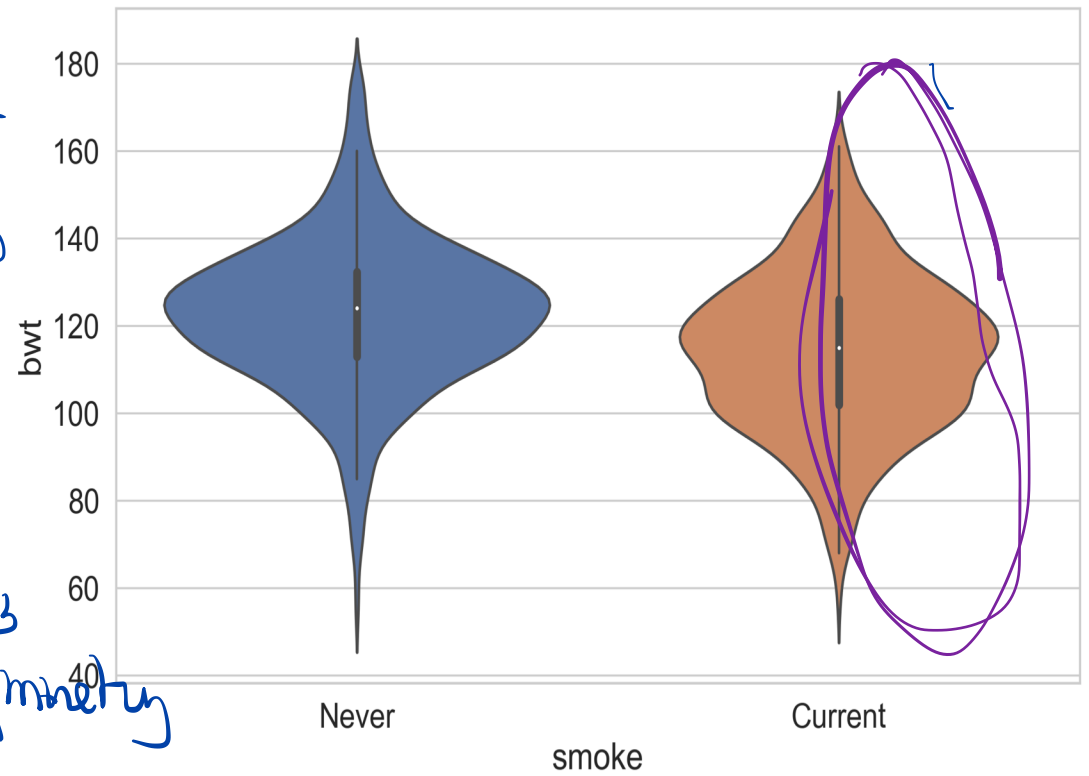Allows us to see shape of the relationship between y and x

# Mix Quantitative & Qualitative

# Mix Quantitative & Qualitative
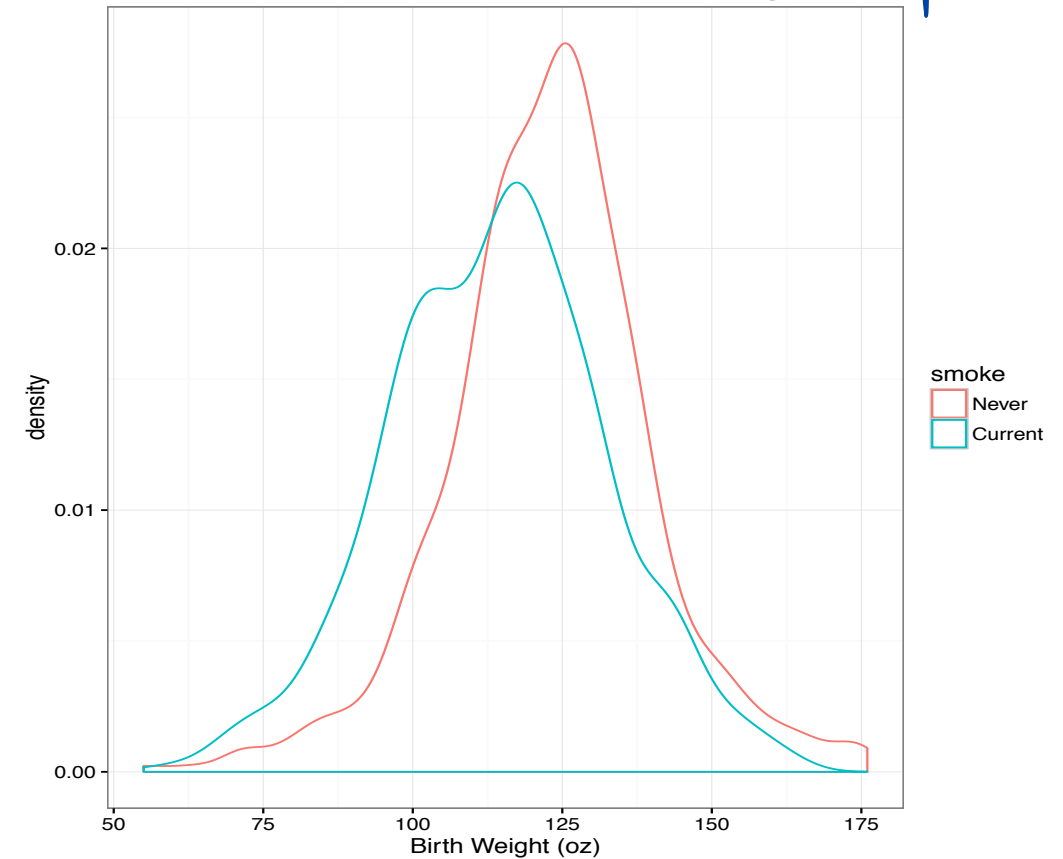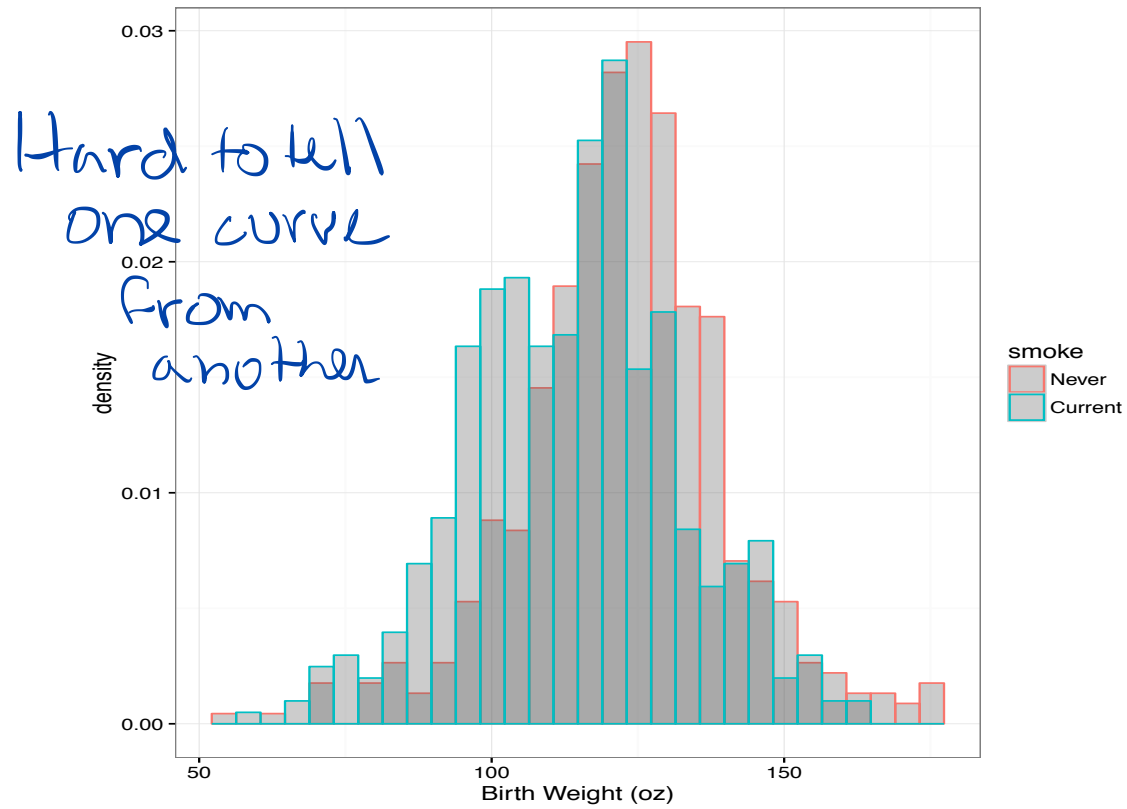
Side-by-side Boxplots



Side-by-side violin plots

Density on its side

Quantitative on y-axis

Can't see modes

Can see tails and assymetry

Can put Qualitative Variable on the x-axis

# Mix of Qualitative and Quantitative

## Overlaid bars/curves

Much preferred

Hard to tell one curve from another

# Two Qualitative Variables

# Pairs of Qualitative Variables
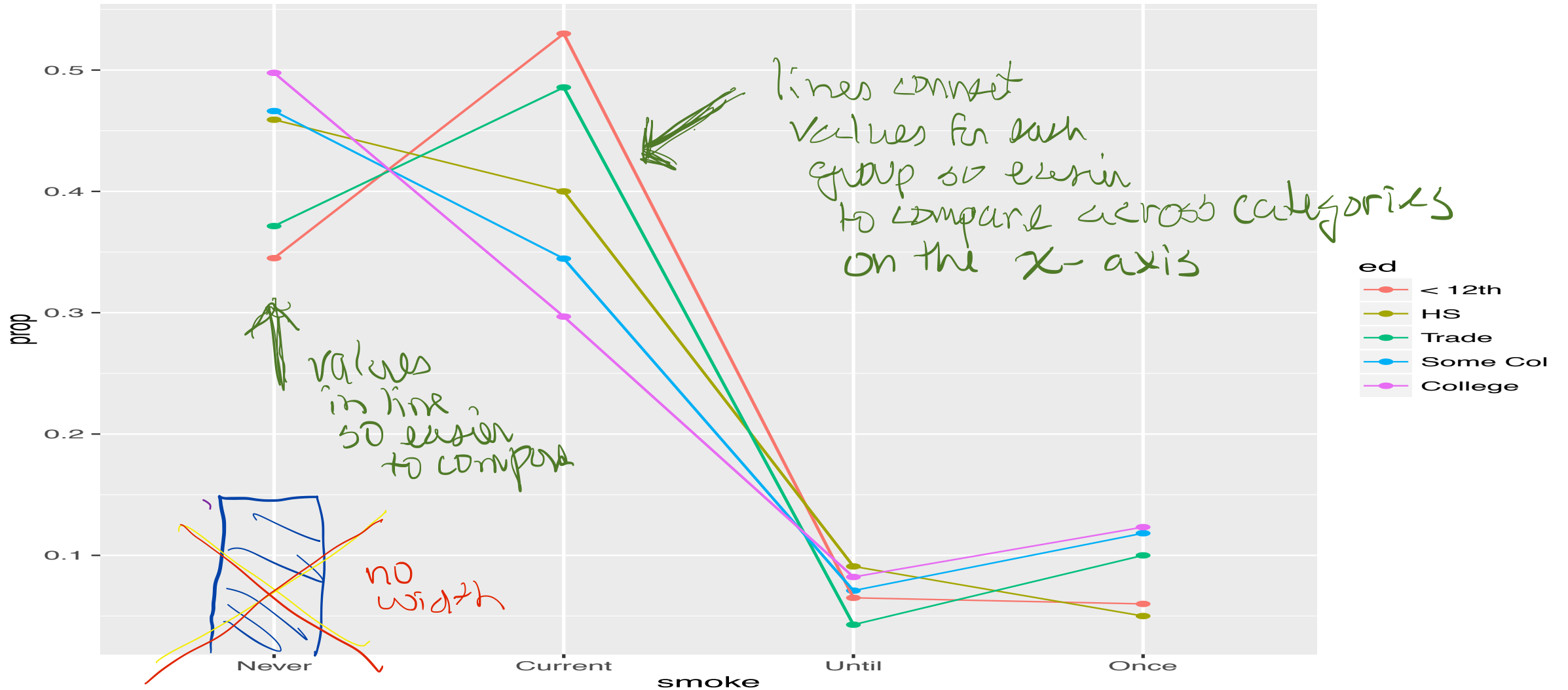


What's the difference between these 2 plots?

# Interaction/Factor Plot

Smoking status normalized within Education level

# Univariate Graphical Displays

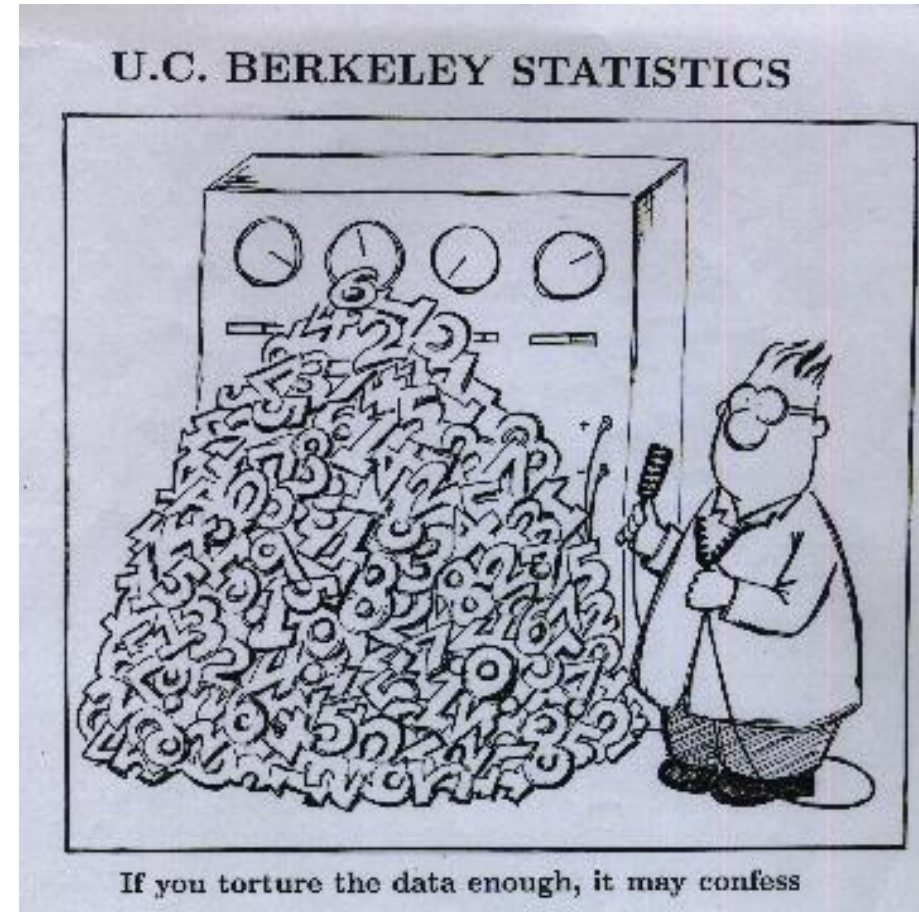| Type | Plot |
|------|------|
| Numeric – | few observations<br>Histogram, Density curve<br>Box plot, Violin plot<br>Normal quantile plot<br>Few Observations - Rug plot, Dot plot<br>Caution if discrete: density curves and box plots may be misleading |
| Categorical –<br>Counts of categories | Dot chart<br>Bar chart<br>Pie chart (avoid!)<br>Caution if ordinal –order of bars, dots, etc. should reflect category order |

# Bivariate Graphical Displays

|  | Numeric | Categorical |
|---|---|---|
| Numeric | Scatter plot<br>Smooth scatter<br>Contour plot<br>Smooth lines and curves | Multiple histograms,<br>density curves,<br>Avoid jiggling! |
| Categorical |  | Side-by-side bar plot<br>Overlaid Lines plot<br>Side-by-side dot chart<br>Mosaic plot<br>Avoid stacking! |

# Caution about EDA

With enough data, if you look hard enough you will find something *"interesting"*

Important to differentiate **inferential conclusions** about world from **exploratory analysis of data**

# Take care with EDA

➢ EDA can provide valuable insights about the data and data collection process

BUT

➢ Be cautious about drawing/reporting conclusions
  ➢ Recognize that EDA biases your view
  ➢ Be careful about sharing plots or hypothesis without additional validation …

➢ Have a lot of data?  Apply EDA to sample of the data before conducting formal analysis.