# Introduction to data science & artificial intelligence (INF7100)

Arthur Charpentier

#211 Statistical Functions (cdf and density)

été 2020

# Statistical Functions



Nicole Oresme, Tractatus de latitudinibus formarum, 1486

# Cumulative Distribution Function

Given a random variable $X$, $F(y) = \mathbb{P}[X \leq x]$

$F$ is an increasing function, taking values in $[0, 1]$.

Consider a sample $\boldsymbol{x} = \{x_1, y_2, \cdots, x_n\}$, a natural estimator is

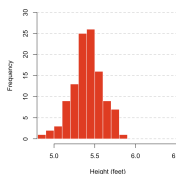$$\widehat{F}(x) = \frac{1}{n} \sum_{i=1}^{n} \mathbf{1}(x_i \leq x)$$

```python
1 > import numpy
2 > x = numpy.sort(x)
3 > n = x.size
4 > y = numpy.arange(1, n+1) / n
5 > import matplotlib.pyplot as plt
6 > plt.plot(x, np.linspace(0, 1, n, endpoint=False))
```
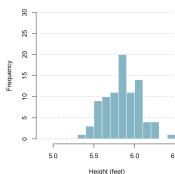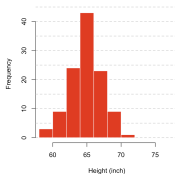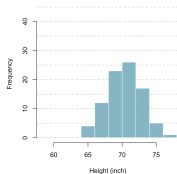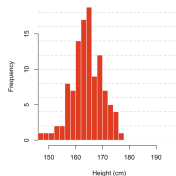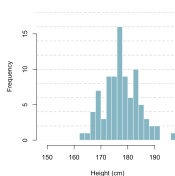
```r
1 > x = sort(x)
2 > n = length(x)
3 > y = (1:n)/n
4 > plot(ecdf(x))
```

# Density & Histogram

Given a random variable $X$, $f$ is such that $F(x) = \displaystyle\int_{-\infty}^{x} f(t)dt$

or conversely, $f(x) = F'(x)$.

Thus, $\mathbb{P}(X \in [a, b]) = \displaystyle\int_{a}^{b} f(t)dt$

# Histogram & Density

Can be used to compare distributions

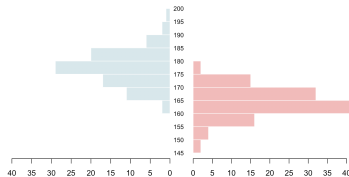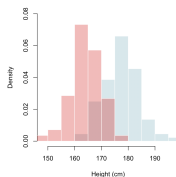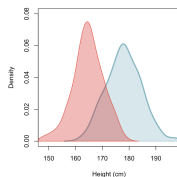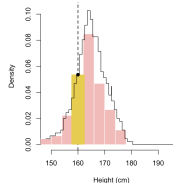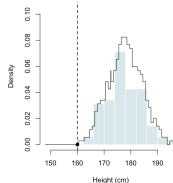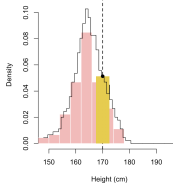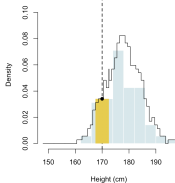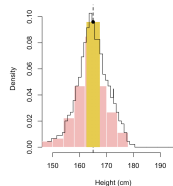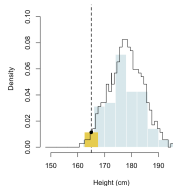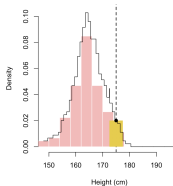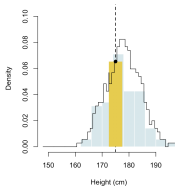# Moving Histogram

$$\widehat{f_h}(x) = \frac{1}{nh} \sum_{i=1}^{n} \mathbf{1}(|x_i - x| \leq h/2)$$

## Moving Histogram

$\widehat{F}$ cannot be differentiated, but we can consider

$$f_h(x) = \frac{1}{h} \underbrace{F(x + h/2) - F(x - h/2)}_{\mathbb{P}(X \in [x \pm h/2])}$$

i.e.

$$\widehat{f_h}(x) = \frac{1}{nh} \sum_{i=1}^{n} \mathbf{1}\left(x_i \in [x - h/2, x + h/2]\right)$$
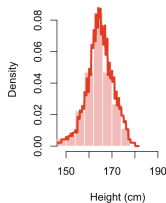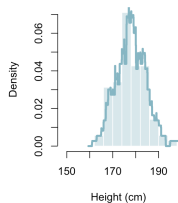
$$\widehat{f_h}(x) = \frac{1}{nh} \sum_{i=1}^{n} \mathbf{1}(|x_i - x| \leq h/2)$$

One can prove that $\mathbb{E}(\widehat{f_h}(x)) = f_h(x) \sim f(x) + \dfrac{h^2}{24} f''(x)$

i.e. bias$(\widehat{f_h}(x)) \sim \dfrac{h^2}{24} f''(x)$, while Var$(\widehat{f_h}(x)) \sim \dfrac{1}{nh} \cdot f_h(x)$
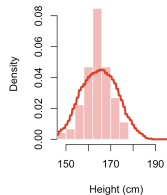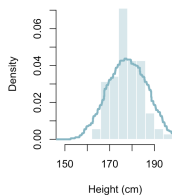
# Moving Histogram

$$\widehat{f}_h(x) = \frac{1}{nh} \sum_{i=1}^{n} \mathbf{1}(|x_i - x| \leq h/2)$$



small $h$
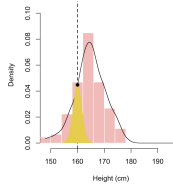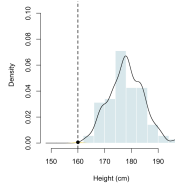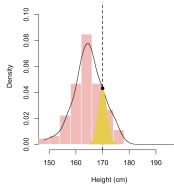bias bias$(\widehat{f}_h(x))$ small
variance Var$(\widehat{f}_h(x))$ large

large $h$
bias bias$(\widehat{f}_h(x))$ large
variance Var$(\widehat{f}_h(x))$ small

# Kernel Density

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

# Kernel Density

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



small $h$
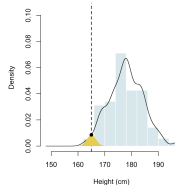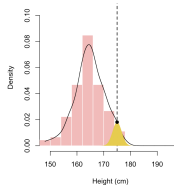bias bias$(\widehat{f}_h(x))$ small
variance Var$(\widehat{f}_h(x))$ large

large $h$
bias bias$(\widehat{f}_h(x))$ large
variance Var$(\widehat{f}_h(x))$ small

# Histogram & Density



```
1 > import matplotlib.pyplot as plt
2 > hist = plt.hist(x, bins=30, normed=True)
3 > from sklearn.neighbors import KernelDensity
4 > k = KernelDensity(bandwidth=1.0, kernel='gaussian')
5 > k.fit(x[:, None])
```
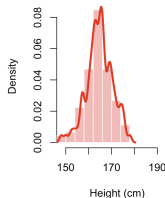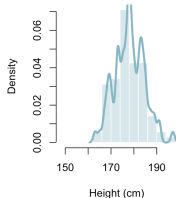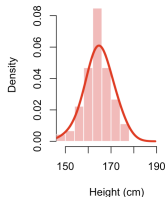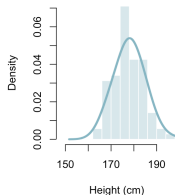
```
1 > hist(x, probability=TRUE)
2 > plot(density(x))
3 > plot(density(x), kernel="gaussian", bw=1)
```

# Histogram & Density



**Distribution of marathon finishing times**

The small spikes are people making their goals, with not a minute to spare. A finishing time of 3:59 is 1.4 times as likely as one of 4:01.
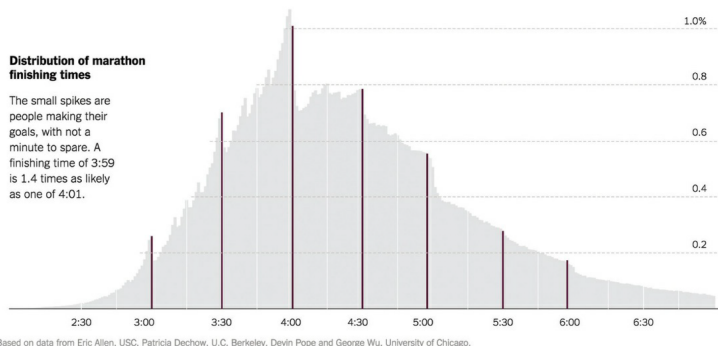
Based on data from Eric Allen, USC, Patricia Dechow, U.C. Berkeley, Devin Pope and George Wu, University of Chicago.

Reference-Dependent Preferences: Evidence from Marathon Runners