

Nonparametric Statistics

1. An inference or model that does not use the probability distribution generating the data
2. Aspects of the probability dist'n may be known, but the complexity of the dist'n is unknown and is adaptive to the data (gets better with more data)

-
- Descriptive statistics and EDA are mostly nonparametric
 - Semiparametric statistical inference: part of the model is parametric, part is nonparametric

Ex: $X_i | \mu_i \sim \text{Normal}(\mu_i, 1) \cdot$
 $\mu_i \sim F$ (arbitrary dist'n) \cdot

- ① Empirical dist'n functions
 - ② Bootstrap
 - ③ Permutation methods — next class meeting
 - ④ Goodness of fit
- method of moments
-

EDFs (empirical distribution functions)

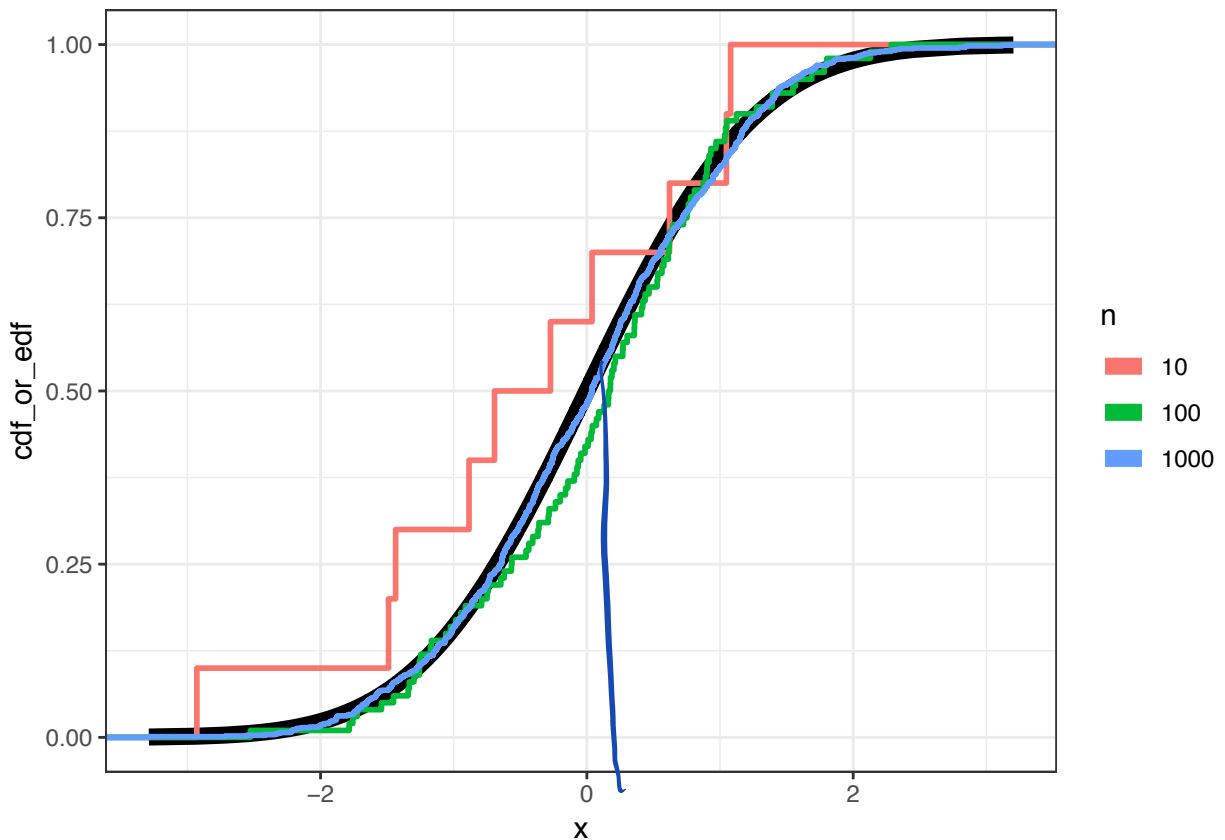
Say $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} F$ (some distribution)

$$\text{Let } 1(X_i \leq y) = \begin{cases} 0 & \text{if } X_i > y \\ 1 & \text{if } X_i \leq y \end{cases}$$

✓ Random variable: $\hat{F}_X(y) = \frac{1}{n} \sum_{i=1}^n 1(X_i \leq y)$

Observed variable: $F_X(y) = \frac{1}{n} \sum_{i=1}^n 1(x_i \leq y)$

EDF from Normal Data



Pointwise Convergence

$$\hat{F}_X(y) \rightarrow F(y) \text{ w/ probability } 1$$

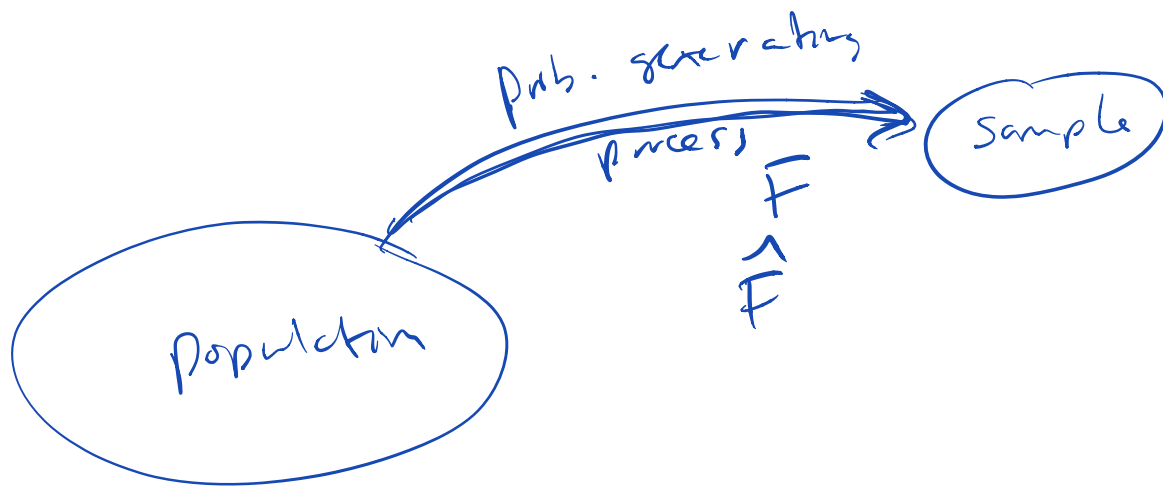
for each y

$n \rightarrow \infty$

Glivenko - Cantelli Theorem

$$\left\{ \sup_{y \in \mathbb{R}} |\hat{F}_X(y) - F(y)| \rightarrow 0 \right.$$

with probability 1



Statistical Functionals

A statistical functional $T(F)$ is any function of cdf, F . Examples:

- $\mu(F) = \int x dF(x)$ (general)
- $= \sum x f(x)$ (discrete)
- $= \int x f(x) dx$ (continuous)

- $\sigma^2(F) = \int (x - \mu(F))^2 dF(x)$

- $m(F) = F^{-1}(1/2)$

Plug-in Estimators of Statistical functionals
from EDFs:

$$\begin{aligned} \bullet \hat{\mu} &= m(\hat{F}) = \int x d\hat{F}(x) \\ &= \sum_{i=1}^n x_i \hat{f}(x_i) \\ &= \sum_{i=1}^n x_i \frac{1}{n} \quad (\text{sample mean}) \end{aligned}$$

$$\bullet \hat{\sigma}^2 = \sigma^2(\hat{F}) = \sum_{i=1}^n (x_i - \hat{\mu})^2 \frac{1}{n}$$

$$\bullet \hat{m} = m(\hat{F}) = \hat{F}^{-1}(1/2)$$

EDF CLT for a statistical functional

$$\left[\frac{T(F) - T(\hat{F})}{\hat{se}(T(\hat{F}))} \right] \xrightarrow{n \rightarrow \infty} \text{Normal}(0, 1)$$

Linear Statistical Functional:

$$T(F) = \int a(x) dF(x)$$

$$\begin{aligned} \text{Var}(T(\hat{F})) &= \frac{1}{n^2} \sum_{i=1}^n \text{Var}(a(x_i)) \\ &= \frac{\text{Var}(a(x))}{n} \end{aligned}$$

$$\text{se}(T(\hat{F})) = \sqrt{\frac{\text{Var}(a(x))}{n}}$$

$$\hat{\text{se}}(T(\hat{F})) = \sqrt{\frac{\text{Var}_{\hat{F}}(a(x))}{n}}$$

Bootstrap

Basic idea: Use \hat{F} (EDF) in place of F to get sampling distributions

Bootstrap Sample

$X_1, X_2, \dots, X_n \stackrel{iid}{\sim} F$

from \hat{F} EPF

If I want n iid observations from \hat{F} , sample n observations with replacement from $\{X_1, X_2, \dots, X_n\}$

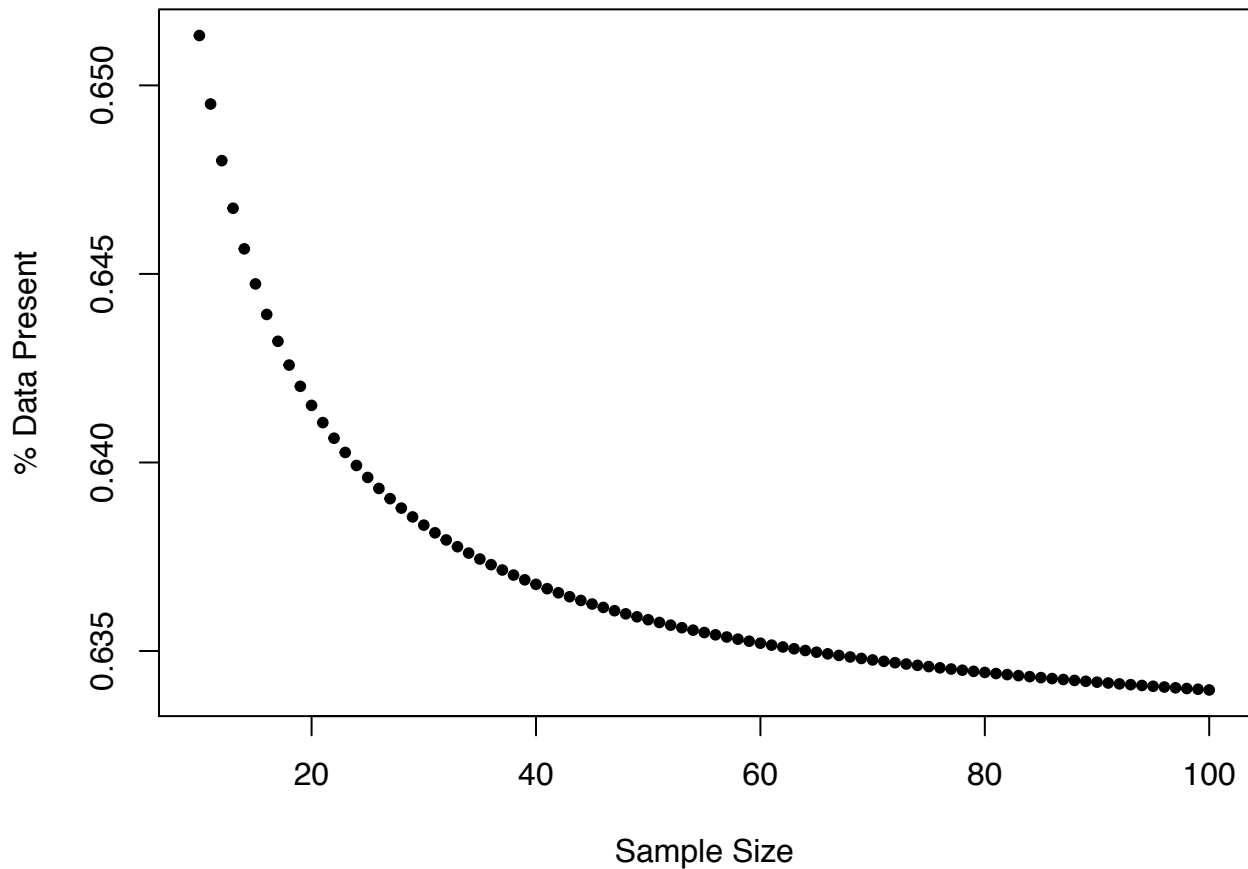
\hat{F} puts probability $1/n$ on each X_i

Probability of not being sampled:

$$\underline{\underline{\left(1 - \frac{1}{n}\right)^n}}$$

Percentage of Data Present in a Bootstrap Sample

For a sample of size n , what percentage of the data is present in any given bootstrap sample?



Suppose we're interested in $\theta = T(F)$.
We estimate it by $\hat{\theta} = T(\hat{F}_x)$

Idea:

For $b = 1, 2, \dots, B$ we draw bootstrap
data sets $x_1^{*(b)}, x_2^{*(b)}, \dots, x_n^{*(b)}$.

Example $n = 100$, $B = 10,000$

We can calculate the estimator $\hat{\theta}$ on each bootstrap sample:

$$\hat{\theta}^{*(1)}, \hat{\theta}^{*(2)}, \dots, \hat{\theta}^{*(B)}$$

Three ways to use bootstrap samples to get confidence intervals:

① Percentile intervals

② Pivotal intervals

③ Studentized pivotal intervals

1-2 Bootstrap ^{two-sided} Confidence Intervals

1st moment
pivotal

2nd moment
pivotal

① Percentile interval:

Let $p_{\alpha/2}^*$ and $p_{1-\alpha/2}^*$ be
the $\alpha/2$ and $1-\alpha/2$ percentiles
of $\hat{\theta}^{(1)}, \hat{\theta}^{(2)}, \dots, \hat{\theta}^{(B)}$.

CI is then:

$$(p_{\alpha/2}^*, p_{1-\alpha/2}^*)$$

Aside Suppose $X_1, X_2, \dots, X_n \sim \text{Normal}(0, 1)$

EDF $\hat{F} \sim \text{Normal}(\bar{X}, S^2)$

residuals

BS $X_i - \bar{X} \Rightarrow \hat{F} \sim \text{Normal}(0, S^2)$

BS $\frac{X_i - \bar{X}}{\sqrt{S^2}} \Rightarrow \hat{F} \sim \text{Normal}(0, 1)$

Studentized
residuals

$Z_{\alpha/2} \quad Z_{1-\alpha/2}$

② Pivotal interval (1st moment pivotal interval):

We calculate percentiles on

$\hat{\theta}^{*(h)} - \hat{\theta}$, call them

q_{α}^* . They are bootstrap

estimates of q_{α} , which are

the α percentiles of $\hat{\theta} - \theta$

If we know q_{α} then the

following is a $1-\alpha$ CI:

$$(\hat{\theta} - q_{1-\alpha/2}, \hat{\theta} - q_{\alpha/2}) \checkmark$$

$$\begin{aligned} 1-\alpha &= \Pr(q_{\alpha/2} \leq \hat{\theta} - \theta \leq q_{1-\alpha/2}) \\ &= \Pr(-q_{1-\alpha/2} \leq \theta - \hat{\theta} \leq -q_{\alpha/2}) \\ &= \Pr(\hat{\theta} - q_{1-\alpha/2} \leq \theta \leq \hat{\theta} - q_{\alpha/2}) \end{aligned}$$

Recall $\hat{\theta}^* - \hat{\theta}$ is approx. to $\hat{\theta} - \theta$

Suppose p_α^* is the α percentile of $\hat{\theta}^*$. Then $p_\alpha^* - \hat{\theta}$ is the approx. α percentile of $\hat{\theta} - \theta$

Therefore $p_\alpha^* - \hat{\theta}$ is the b.s. estimate of q_α . Plugging this into the above, we get the $(1-\alpha)$ CI is:

$$(2\hat{\theta} - p_{1-\alpha/2}^*, 2\hat{\theta} - p_{\alpha/2}^*) \checkmark$$

$$= (\hat{\theta} - q_{1-\alpha/2}^*, \hat{\theta} - q_{\alpha/2}^*) \checkmark$$

③ Studentized pivotal intervals (2nd moment pivotal)

The goal is to approximate the sampling dist'n of

$$\frac{\hat{\theta} - \theta}{\text{se}(\hat{\theta})}$$

Approximated by:

$$\frac{\hat{\theta}^* - \hat{\theta}}{\text{se}(\hat{\theta}^*)}$$

Let z_{α}^* be the α percentile of

$$\left\{ \frac{\hat{\theta}^{*(1)} - \hat{\theta}}{\text{se}(\hat{\theta}^{*(1)})}, \frac{\hat{\theta}^{*(2)} - \hat{\theta}}{\text{se}(\hat{\theta}^{*(2)})}, \dots, \frac{\hat{\theta}^{*(B)} - \hat{\theta}}{\text{se}(\hat{\theta}^{*(B)})} \right\}$$

Example: $\hat{\theta} = \bar{x}$
 $\hat{\theta}^* = \bar{x}^*$
 $\hat{se}(\hat{\theta}^*) = \frac{s^*}{\sqrt{n}}$

$$\frac{\hat{\theta}^* - \hat{\theta}}{\hat{se}(\hat{\theta}^*)} = \frac{\bar{x}^* - \bar{x}}{s^*/\sqrt{n}}$$

The $(1-\alpha)$ two-sided b.s. CI is

$$\left(\hat{\theta} - z_{1-\alpha/2}^* \hat{se}(\hat{\theta}), \hat{\theta} + z_{\alpha/2}^* \hat{se}(\hat{\theta}) \right)$$



replacing Normal(0,1)
percentiles

How do we get $\hat{se}(\hat{\theta})$ in
nonparametric settings?

$$\hat{se}(\hat{\theta}) = \sqrt{\frac{1}{B} \sum_{b=1}^B (\hat{\theta}^{*(b)} - \frac{1}{B} \sum_{k=1}^B \hat{\theta}^{*(k)})^2}$$

But how to get $\hat{se}(\hat{\theta}^{*(b)})$???

$$\text{Exp}(\lambda) \quad E[X] = \frac{1}{\lambda} = 2$$

$$\lambda = 1/2$$

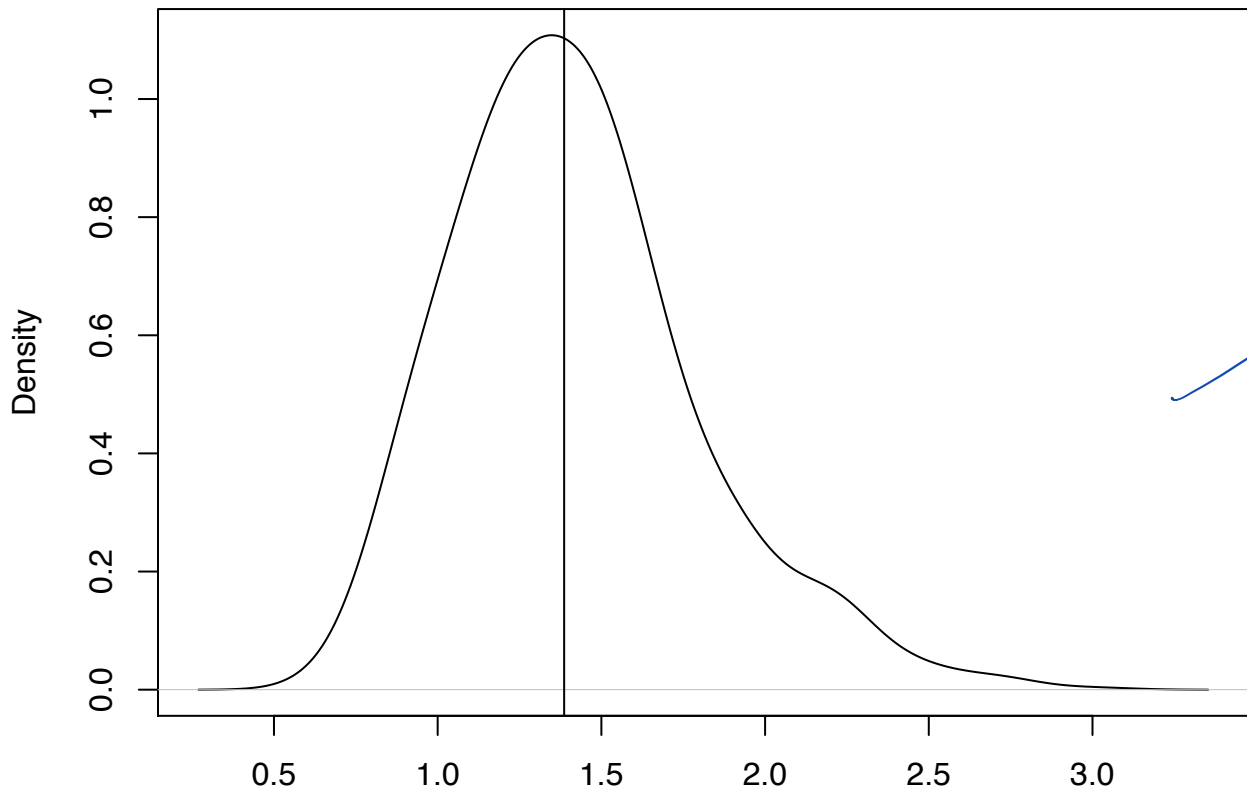
Example: Bootstrap on Exponential Data

In the homework, you will be performing a bootstrap t-test of the mean and a bootstrap percentile CI of the median for the following Exponential(λ) data:

```
> set.seed(1111)
> pop_mean <- 2
> X <- matrix(rexp(1000*30, rate=1/pop_mean), nrow=1000, ncol=30)
```

Let's construct a pivotal bootstrap CI of the median here instead.

```
> # population median 2*log(2)
> pop_med <- qexp(0.5, rate=1/pop_mean); pop_med
[1] 1.386294
>
> obs_meds <- apply(X, 1, median)
> plot(density(obs_meds, adj=1.5), main=" "); abline(v=pop_med)
```



N = 1000 Bandwidth = 0.1155

Some embarrassingly inefficient code to calculate bootstrap medians.

```
> B <- 1000
> boot_meds <- matrix(0, nrow=1000, ncol=B)
>
> for(b in 1:B) {
```



```

+   idx <- sample(1:30, replace=TRUE)
+   boot_meds[,b] <- apply(X[,idx], 1, median)
+ }

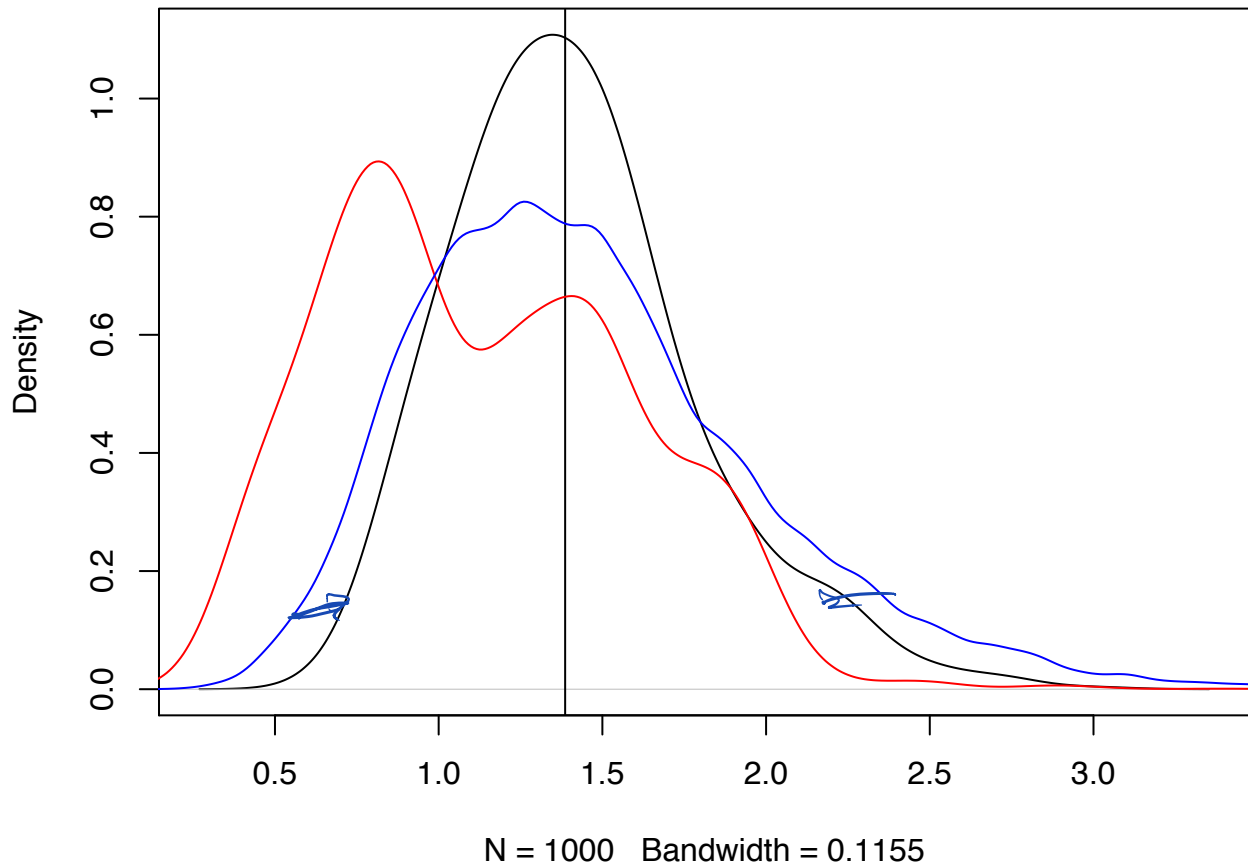
```

Plot the bootstrap medians.

```

> plot(density(obs_meds, adj=1.5), main=" "); abline(v=pop_med)
> lines(density(as.vector(boot_meds[1:4,]), adj=1.5), col="red")
> lines(density(as.vector(boot_meds), adj=1.5), col="blue")

```



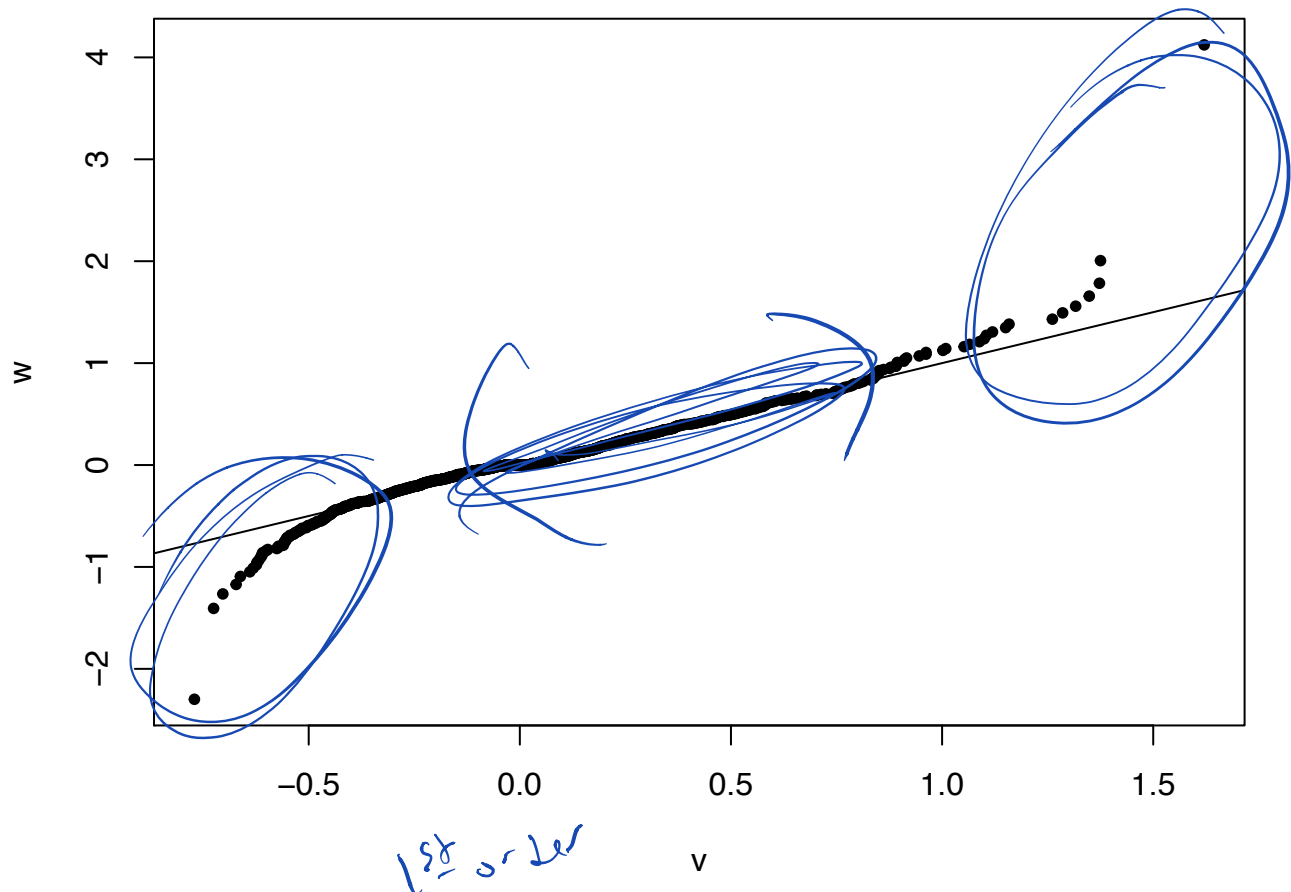
Compare sampling distribution of $\hat{\theta} - \theta$ to $\hat{\theta}^* - \hat{\theta}$.

```

> v <- obs_meds - pop_med
> w <- as.vector(boot_meds - obs_meds)
> qqplot(v, w, pch=20); abline(0,1)

```

Handwritten blue notes: $\hat{\theta} - \theta$ and $\hat{\theta}^* - \hat{\theta}$ with arrows pointing to the corresponding terms in the code above.



Does a 95% bootstrap pivotal interval provide coverage?

```

> ci_lower <- apply(boot_meds, 1, quantile, probs=0.975)
> ci_upper <- apply(boot_meds, 1, quantile, probs=0.025)
>
> ci_lower <- 2*obs_meds - ci_lower
> ci_upper <- 2*obs_meds - ci_upper
>
> ci_lower[1]; ci_upper[1]
[1] 0.8958224
[1] 2.113859
>
> cover <- (pop_med >= ci_lower) & (pop_med <= ci_upper)
> mean(cover)
[1] 0.809
>
> # :-(

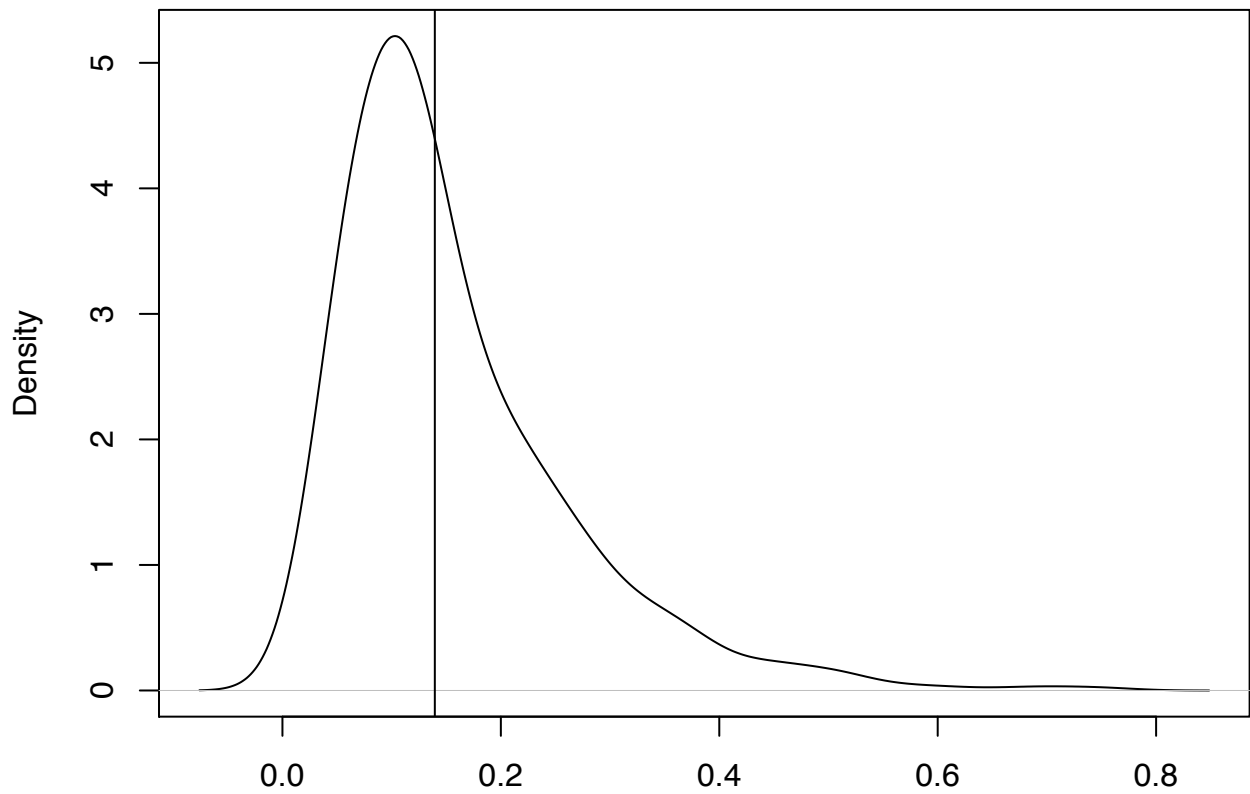
```

Let's check the bootstrap variances.

```

> sampling_var <- var(obs_meds)
> boot_var <- apply(boot_meds, 1, var)
> plot(density(boot_var, adj=1.5), main=" ")
> abline(v=sampling_var)

```



N = 1000 Bandwidth = 0.0303

We repeated this simulation over a range of n and B .

n	B	coverage	avg CI width
1e+02	1000	0.868	0.7805404
1e+02	2000	0.872	0.7882278
1e+02	4000	0.865	0.7852837
1e+02	8000	0.883	0.7817222
1e+03	1000	0.923	0.2465840
1e+03	2000	0.909	0.2477463
1e+03	4000	0.915	0.2475550
1e+03	8000	0.923	0.2458167
1e+04	1000	0.935	0.0781421
1e+04	2000	0.937	0.0784541
1e+04	4000	0.942	0.0784559
1e+04	8000	0.948	0.0785591
1e+05	1000	0.949	0.0246918
1e+05	2000	0.942	0.0246938

Goodness of Fit Methods

We don't know the dist'n of the data,
but we'd like to test or assess
its fit to a known distribution

- ① Chi-square Gof
- ② KS Test
- ③ Method of moments

Chi-square Gof

$X_1, X_2, \dots, X_n \stackrel{iid}{\sim} F$

Test: $H_0: F \in \{F_\theta: \theta \in \Theta\}$

$H_1: \text{not } H_0$

Divide the support of $\{F_\theta: \theta \in \Theta\}$
in K bins I_1, I_2, \dots, I_K

Example : Normal (μ, σ^2)

$$\begin{array}{ccccccc} (-\infty, -10), & (-10, -9) & \dots & (9, 10), & (10, \infty) \\ I_1 & I_2 & \dots & I_{k-1} & I_k \end{array}$$

For $j = 1, 2, \dots, k$ calculate

$$q_j(\theta) = \int_{I_j} dF_\theta(x)$$

Suppose observe data x_1, x_2, \dots, x_n .

Let n_j be the number of data points in interval I_j .

Let $\tilde{\theta}$ be the value of θ

that is the MLE of :

$$\prod_{j=1}^k q_j(\theta)^{n_j}$$

Form GOF statistic:

$$S(x) = \sum_{j=1}^K \frac{(n_j - n q_j(\bar{\theta}))^2}{n q_j(\bar{\theta})}$$

$n q_j(\bar{\theta})$ is the expected number of observations in I_j with parameter values $\bar{\theta}$

When H_0 is true, $S(x)$ has

a χ^2_v where $v = K - \dim(\theta)$

p-value = $\Pr(S(x^*) \geq s(x))$

where $S(x^*) \sim \chi^2_v$.

Goodness of Fit Example: Hardy-Weinberg

Suppose at your favorite SNP, we observe genotypes from 100 randomly sampled individuals as follows:

AA	AT	TT
28	60	12

If we code these genotypes as 0, 1, 2, testing for Hardy-Weinberg equilibrium is equivalent to testing whether $X_1, X_2, \dots, X_{100} \stackrel{\text{iid}}{\sim} \text{Binomial}(2, \theta)$ for some unknown allele frequency of T, θ .

The parameter dimension is such that $d = 1$. We will also set $k = 3$, where each bin is a genotype. Therefore, we have $n_1 = 28$, $n_2 = 60$, and $n_3 = 12$. Also,

$$q_1(\theta) = (1 - \theta)^2, \quad q_2(\theta) = 2\theta(1 - \theta), \quad q_3(\theta) = \theta^2.$$

Forming the multinomial likelihood under these bin probabilities, we find $\tilde{\theta} = (n_2 + 2n_3)/(2n)$. The degrees of freedom of the χ_v^2 null distribution is $v = k - d - 1 = 3 - 1 - 1 = 1$.

Let's carry out the test in R.

```
> n <- 100
> nj <- c(28, 60, 12)
>
> # parameter estimates
> theta <- (nj[2] + 2*nj[3])/(2*n)
> qj <- c((1-theta)^2, 2*theta*(1-theta), theta^2)
>
> # gof statistic
> s <- sum((nj - n*qj)^2 / (n*qj))
>
> # p-value
> 1-pchisq(s, df=1)
[1] 0.02059811
```

Kolmogorov-Smirnov Test

- ① Form EDF \hat{F}
- ② Parametric F_θ (θ known)
- ③ Form statistic:

$$D(X) = \max_y \left| \hat{F}_X(y) - F_\theta(y) \right|$$

Null distribution of $D(X)$ is known, based on Brownian bridge.

$$H_0: F = F_\theta \text{ vs. } H_1: F \neq F_\theta$$

Two-sample KS-test

$$X_1, \dots, X_n \sim F_X$$

$$Y_1, \dots, Y_m \sim F_Y$$

$$H_0: F_X = F_Y \text{ vs. } H_1: F_X \neq F_Y$$

$$D(x, y) = \max_z |\hat{F}_x(z) - \hat{F}_y(z)|$$

When H_0 is true, one can
calculate the dist'n of

$$D(x, y) .$$

KS Test Example: Exponential vs Normal

```
ks.test(x, y, ...,  
        alternative = c("two.sided", "less", "greater"),  
        exact = NULL)
```

Two sample KS test.

```
> x <- rnorm(100, mean=1)  
> y <- rexp(100, rate=1)  
> wilcox.test(x, y)
```

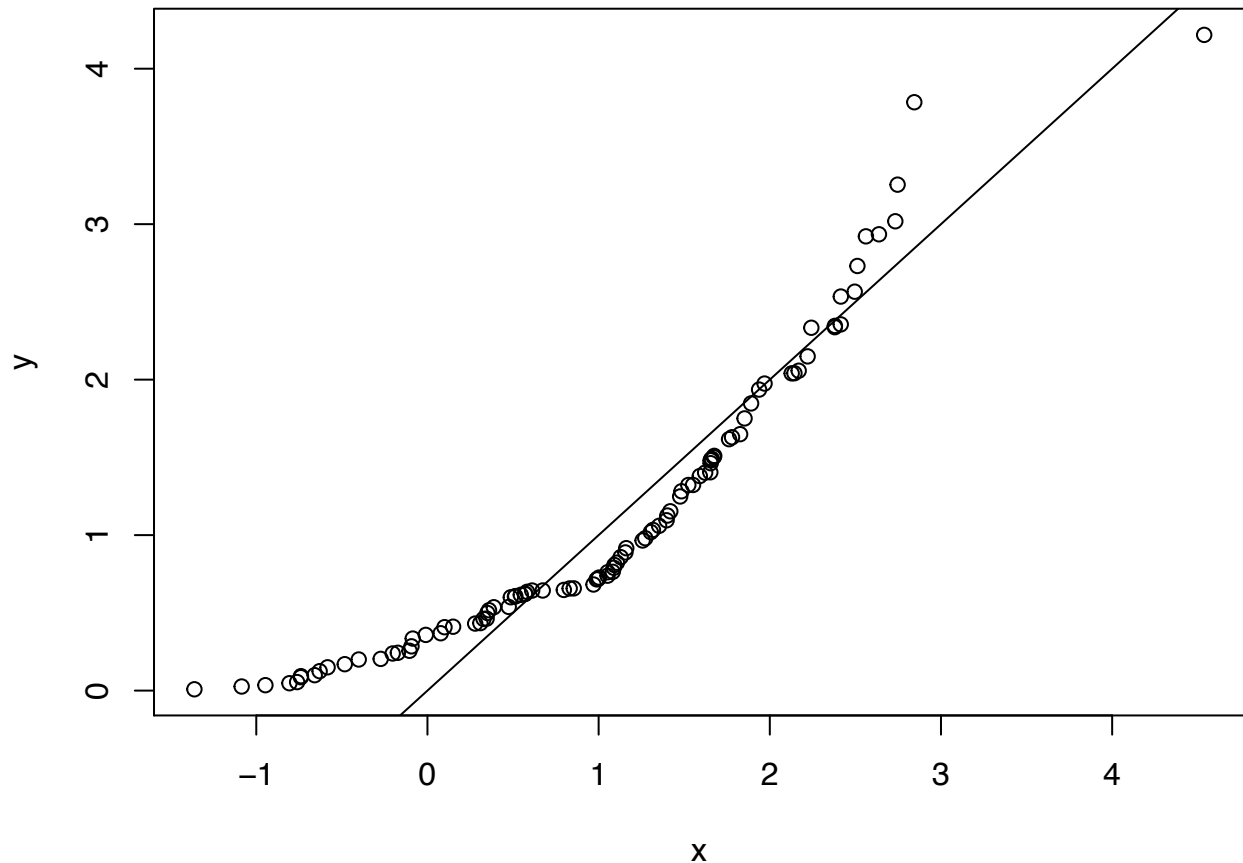
Wilcoxon rank sum test with continuity correction

```
data: x and y  
W = 5021, p-value = 0.9601  
alternative hypothesis: true location shift is not equal to 0  
> ks.test(x, y)
```

Two-sample Kolmogorov-Smirnov test

```
data: x and y  
D = 0.19, p-value = 0.0541  
alternative hypothesis: two-sided
```

```
> qqplot(x, y); abline(0,1)
```



One sample KS tests.

```
> ks.test(x=x, y="pnorm")
```

One-sample Kolmogorov-Smirnov test

```
data: x
D = 0.41398, p-value = 2.554e-15
alternative hypothesis: two-sided
```

```
> ks.test(x=x, y="pnorm", mean=1)
```

One-sample Kolmogorov-Smirnov test

```
data: x
D = 0.068035, p-value = 0.7436
alternative hypothesis: two-sided
```

Standardize (mean center, sd scale) the observations before comparing to a Normal(0,1) distribution.

```
> ks.test(x=((x-mean(x))/sd(x)), y="pnorm")
```

One-sample Kolmogorov-Smirnov test

```
data: ((x - mean(x))/sd(x))  
D = 0.05896, p-value = 0.8778  
alternative hypothesis: two-sided  
>  
> ks.test(x=((y-mean(y))/sd(y)), y="pnorm")
```

One-sample Kolmogorov-Smirnov test

```
data: ((y - mean(y))/sd(y))  
D = 0.14439, p-value = 0.03092  
alternative hypothesis: two-sided
```