

# Statistical Modeling

Larry Wasserman  
Department of Statistics and Data Science  
and  
Machine Learning Department  
Carnegie Mellon University

# Outline

# Outline

## 1. Parametric Models

# Outline

1. Parametric Models
2. Misspecified Parametric Models

# Outline

1. Parametric Models
2. Misspecified Parametric Models
3. Nonparametric Models

# Outline

1. Parametric Models
2. Misspecified Parametric Models
3. Nonparametric Models
4. Universal Inference (time permitting)

# Parametric Models

## Parametric Models

$$X_1, \dots, X_n \sim p(x; \theta)$$



# Parametric Models

$$X_1, \dots, X_n \sim p(x; \theta)$$

Likelihood:

$$\mathcal{L}(\theta) = \prod_{i=1}^n p(X_i; \theta)$$

# Parametric Models

$$X_1, \dots, X_n \sim p(x; \theta)$$

Likelihood:

$$\mathcal{L}(\theta) = \prod_{i=1}^n p(X_i; \theta)$$

Maximum likelihood (MLE):

# Parametric Models

$$X_1, \dots, X_n \sim p(x; \theta)$$

Likelihood:

$$\mathcal{L}(\theta) = \prod_{i=1}^n p(X_i; \theta)$$

Maximum likelihood (MLE):

Find  $\hat{\theta}$  to maximize  $\mathcal{L}(\theta)$

## Parametric Models

$$X_1, \dots, X_n \sim p(x; \theta)$$

Likelihood:

$$\mathcal{L}(\theta) = \prod_{i=1}^n p(X_i; \theta)$$

Maximum likelihood (MLE):

Find  $\hat{\theta}$  to maximize  $\mathcal{L}(\theta)$

$$\hat{\theta} \approx N(\theta, J)$$

where  $J = I^{-1}$  and  $I$  is the Fisher information:

$$I_{jk} = -\mathbb{E} \left[ \frac{\partial^2 \log \mathcal{L}(\theta)}{\partial \theta_j \partial \theta_k} \right].$$

# Parametric Models

## Parametric Models

Let

$$C = \left\{ \theta : \frac{\mathcal{L}(\theta)}{\mathcal{L}(\hat{\theta})} > e^{-(1/2)\chi_{d,\alpha}^2} \right\}.$$

## Parametric Models

Let

$$C = \left\{ \theta : \frac{\mathcal{L}(\theta)}{\mathcal{L}(\hat{\theta})} > e^{-(1/2)\chi_{d,\alpha}^2} \right\}.$$

Then

$$P(\theta \in C) \rightarrow 1 - \alpha.$$

## Parametric Models

Let

$$C = \left\{ \theta : \frac{\mathcal{L}(\theta)}{\mathcal{L}(\hat{\theta})} > e^{-(1/2)\chi_{d,\alpha}^2} \right\}.$$

Then

$$P(\theta \in C) \rightarrow 1 - \alpha.$$

For  $\theta_j$  use  $\hat{\theta}_j \pm z_{\alpha/2} \sqrt{J_{jj}}$ .



## Parametric Models

Let

$$C = \left\{ \theta : \frac{\mathcal{L}(\theta)}{\mathcal{L}(\hat{\theta})} > e^{-(1/2)\chi_{d,\alpha}^2} \right\}.$$

Then

$$P(\theta \in C) \rightarrow 1 - \alpha.$$

For  $\theta_j$  use  $\hat{\theta}_j \pm z_{\alpha/2} \sqrt{J_{jj}}$ .

Or use the profile likelihood:

## Parametric Models

Let

$$C = \left\{ \theta : \frac{\mathcal{L}(\theta)}{\mathcal{L}(\hat{\theta})} > e^{-(1/2)\chi_{d,\alpha}^2} \right\}.$$

Then

$$P(\theta \in C) \rightarrow 1 - \alpha.$$

For  $\theta_j$  use  $\hat{\theta}_j \pm z_{\alpha/2} \sqrt{J_{jj}}$ .

Or use the profile likelihood:

$$\mathcal{L}(\theta_j) = \sup_{\theta_{-j}} \mathcal{L}(\theta)$$

## MLE: some useful facts

## MLE: some useful facts

In general, the mle of  $\psi(\theta)$  is  $\psi(\hat{\theta})$ .

## MLE: some useful facts

In general, the mle of  $\psi(\theta)$  is  $\psi(\hat{\theta})$ .

The MLE is asymptotically linear

## MLE: some useful facts

In general, the mle of  $\psi(\theta)$  is  $\psi(\hat{\theta})$ .

The MLE is asymptotically linear

$$\sqrt{n}(\psi(\hat{\theta}) - \psi(\theta)) = \frac{1}{\sqrt{n}} \sum_i \varphi(X_i) + o_P(1) \rightsquigarrow N(0, \mathbb{E}[\varphi(Z)\varphi(Z)^T])$$

where

$$\varphi(z) = \phi'(\theta) I_{\theta}^{-1} \ell'_{\theta}(x)$$

is the influence function and  $\ell_{\theta}(x) = \log p_{\theta}(x)$ .

## MLE: some useful facts

## MLE: some useful facts

The MLE is optimal:



## MLE: some useful facts

The MLE is optimal:

if  $\tilde{\psi}$  is another (regular) estimator then

## MLE: some useful facts

The MLE is optimal:

if  $\tilde{\psi}$  is another (regular) estimator then

$$\sqrt{n}(\tilde{\psi} - \psi(\hat{\theta})) \rightsquigarrow N(0, \mathbb{E}[\varphi(Z)\varphi(Z)^T]) + \text{noise}$$

## MLE: some useful facts

The MLE is optimal:

if  $\tilde{\psi}$  is another (regular) estimator then

$$\sqrt{n}(\tilde{\psi} - \psi(\hat{\theta})) \rightsquigarrow N(0, \mathbb{E}[\varphi(Z)\varphi(Z)^T]) + \text{noise}$$

But this assumes:

## MLE: some useful facts

The MLE is optimal:

if  $\tilde{\psi}$  is another (regular) estimator then

$$\sqrt{n}(\tilde{\psi} - \psi(\hat{\theta})) \rightsquigarrow N(0, \mathbb{E}[\varphi(Z)\varphi(Z)^T]) + \text{noise}$$

But this assumes:

(1) The model is correct.

## MLE: some useful facts

The MLE is optimal:

if  $\tilde{\psi}$  is another (regular) estimator then

$$\sqrt{n}(\tilde{\psi} - \psi(\hat{\theta})) \rightsquigarrow N(0, \mathbb{E}[\varphi(Z)\varphi(Z)^T]) + \text{noise}$$

But this assumes:

- (1) The model is correct.
- (2) Regularity conditions.

# Examples

Easy:

## Examples

Easy:

$$X_1, \dots, X_n \sim N(\mu, \sigma^2)$$

## Examples

Easy:

$$X_1, \dots, X_n \sim N(\mu, \sigma^2)$$

$$\hat{\mu} = \frac{1}{n} \sum_i X_i$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_i (X_i - \bar{X})^2$$



## Examples

Easy:

$$X_1, \dots, X_n \sim N(\mu, \sigma^2)$$

$$\hat{\mu} = \frac{1}{n} \sum_i X_i$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_i (X_i - \bar{X})^2$$

Hard:

$$p(x) = (1 - \lambda) \underbrace{b(x; \gamma)}_{\text{background}} + \lambda \underbrace{s(x; \theta)}_{\text{signal}}$$

$\hat{\lambda}, \hat{\gamma}, \hat{\theta}$  must be found numerically

Problems: regularity conditions fail. More on this later.

# Problems with the MLE

# Problems with the MLE

Misspecified model

# Problems with the MLE

Misspecified model

Not robust to outliers

# Problems with the MLE

Misspecified model

Not robust to outliers

If there are many nuisance parameters, the MLE will be biased.  
(semiparametric models).

# Problems with the MLE

Misspecified model

Not robust to outliers

If there are many nuisance parameters, the MLE will be biased.  
(semiparametric models).

High-dimensions: MLE no longer optimal.

# Problems with the MLE

Misspecified model

Not robust to outliers

If there are many nuisance parameters, the MLE will be biased.  
(semiparametric models).

High-dimensions: MLE no longer optimal.

Failure of regularity conditions.

# Misspecified Models



# Misspecified Models

Suppose the model is wrong. (It usually is.) What are we estimating?

## Misspecified Models

Suppose the model is wrong. (It usually is.) What are we estimating?

Define the Kullback-Leibler distance

$$K(p, q) = \int p(x) \log \left( \frac{p(x)}{q(x)} \right) dx.$$

The mle converges to  $\theta_*$  where  $\theta_*$  minimizes  $K(p_{\theta}, p)$ .

## Misspecified Models

Suppose the model is wrong. (It usually is.) What are we estimating?

Define the Kullback-Leibler distance

$$K(p, q) = \int p(x) \log \left( \frac{p(x)}{q(x)} \right) dx.$$

The mle converges to  $\theta_*$  where  $\theta_*$  minimizes  $K(p_{\theta}, p)$ .

Is this a good thing?

## Misspecified Models

Suppose the model is wrong. (It usually is.) What are we estimating?

Define the Kullback-Leibler distance

$$K(p, q) = \int p(x) \log \left( \frac{p(x)}{q(x)} \right) dx.$$

The mle converges to  $\theta_*$  where  $\theta_*$  minimizes  $K(p_{\theta}, p)$ .

Is this a good thing?

Not necessarily.

# Misspecified Models

## Misspecified Models

Suppose use the Normal model

$$X_1, \dots, X_n \sim N(\theta, 1).$$

## Misspecified Models

Suppose use the Normal model

$$X_1, \dots, X_n \sim N(\theta, 1).$$

Suppose the true distribution is

$$(1 - \epsilon)N(\theta, 1) + \epsilon\delta_a$$

where  $\delta_a$  is some distribution concentrated at a point  $a$ .

## Misspecified Models

Suppose use the Normal model

$$X_1, \dots, X_n \sim N(\theta, 1).$$

Suppose the true distribution is

$$(1 - \epsilon)N(\theta, 1) + \epsilon\delta_a$$

where  $\delta_a$  is some distribution concentrated at a point  $a$ .

Then  $\theta_* \rightarrow \infty$  as  $a \rightarrow \infty$ .



## Robustness to Model Misspecification

Maximum likelihood (Bayes) are not robust to model misspecification.

## Robustness to Model Misspecification

Maximum likelihood (Bayes) are not robust to model misspecification.

One solution: use **power divergence** (Basu, Harris, Hjort and Jones 1998):

## Robustness to Model Misspecification

Maximum likelihood (Bayes) are not robust to model misspecification.

One solution: use **power divergence** (Basu, Harris, Hjort and Jones 1998):

Minimize

$$S(\theta) = \int p_0(x; \theta)^{1+\alpha} - \left(1 + \frac{1}{\alpha}\right) \frac{1}{n} \sum_i p_0(X_i; \theta)^\alpha$$

## Robustness to Model Misspecification

Maximum likelihood (Bayes) are not robust to model misspecification.

One solution: use **power divergence** (Basu, Harris, Hjort and Jones 1998):

Minimize

$$S(\theta) = \int p_0(x; \theta)^{1+\alpha} - \left(1 + \frac{1}{\alpha}\right) \frac{1}{n} \sum_i p_0(X_i; \theta)^\alpha$$

$\alpha \rightarrow 0$  gives mle.

## Robustness to Model Misspecification

Maximum likelihood (Bayes) are not robust to model misspecification.

One solution: use **power divergence** (Basu, Harris, Hjort and Jones 1998):

Minimize

$$S(\theta) = \int p_0(x; \theta)^{1+\alpha} - \left(1 + \frac{1}{\alpha}\right) \frac{1}{n} \sum_i p_0(X_i; \theta)^\alpha$$

$\alpha \rightarrow 0$  gives mle.

$\alpha > 0$  more robust.

## Robustness to Model Misspecification

Maximum likelihood (Bayes) are not robust to model misspecification.

One solution: use **power divergence** (Basu, Harris, Hjort and Jones 1998):

Minimize

$$S(\theta) = \int p_0(x; \theta)^{1+\alpha} - \left(1 + \frac{1}{\alpha}\right) \frac{1}{n} \sum_i p_0(X_i; \theta)^\alpha$$

$\alpha \rightarrow 0$  gives mle.

$\alpha > 0$  more robust.

$\alpha = 1$  is  $L_2$ :  $\int (p - p_\theta)^2$ .

## Robustness to Model Misspecification

Maximum likelihood (Bayes) are not robust to model misspecification.

One solution: use **power divergence** (Basu, Harris, Hjort and Jones 1998):

Minimize

$$S(\theta) = \int p_0(x; \theta)^{1+\alpha} - \left(1 + \frac{1}{\alpha}\right) \frac{1}{n} \sum_i p_0(X_i; \theta)^\alpha$$

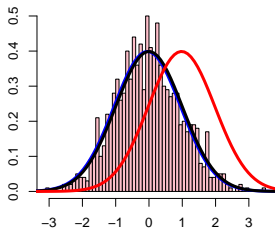
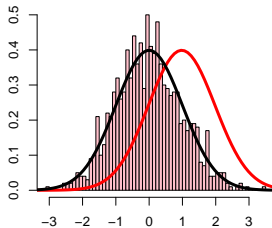
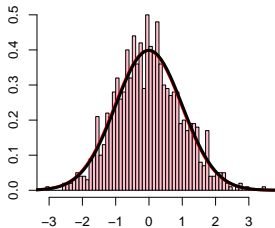
$\alpha \rightarrow 0$  gives mle.

$\alpha > 0$  more robust.

$\alpha = 1$  is  $L_2$ :  $\int (p - p_\theta)^2$ .

$\alpha \uparrow$ : robustness  $\uparrow$  and efficiency  $\downarrow$

# Example





## Simulated Likelihoods

Suppose we want the likelihood ratio

$$R = \prod_{i=1}^n \frac{p(X_i; \theta_1)}{p(X_i; \theta_2)} = \prod_i R(X_i).$$

We may not have a closed form expression It is now common to estimate  $R$  by simulation + classification

## Simulated Likelihoods

$$X_1, \dots, X_N \sim p(x; \theta_1)$$

$$X_{N+1}, \dots, X_{2N} \sim p(x; \theta_2)$$

## Simulated Likelihoods

$$X_1, \dots, X_N \sim p(x; \theta_1)$$

$$X_{N+1}, \dots, X_{2N} \sim p(x; \theta_2)$$

S		1	1	...	1	0	0	...	0
X		$X_1$	$X_2$	...	$X_N$	$X_{N+1}$	$X_{N+2}$	...	$X_{2N}$

# Simulated Likelihoods

## Simulated Likelihoods

Classifier (neural net, random forest, ...)

$$h(x) = P(S = 1|X = x) = \frac{p(x; \theta_1)}{p(x; \theta_1) + p(x; \theta_2)} = \frac{1}{1 + R(x)}$$

## Simulated Likelihoods

Classifier (neural net, random forest, ...)

$$h(x) = P(S = 1|X = x) = \frac{p(x; \theta_1)}{p(x; \theta_1) + p(x; \theta_2)} = \frac{1}{1 + R(x)}$$

$$R(x) = \frac{p(x; \theta_1)}{p(x; \theta_2)} = \frac{h(x)}{1 - h(x)}$$

## Simulated Likelihoods

Classifier (neural net, random forest, ...)

$$h(x) = P(S = 1|X = x) = \frac{p(x; \theta_1)}{p(x; \theta_1) + p(x; \theta_2)} = \frac{1}{1 + R(x)}$$

$$R(x) = \frac{p(x; \theta_1)}{p(x; \theta_2)} = \frac{h(x)}{1 - h(x)}$$

$$\hat{R} = \prod_{i=1}^n \frac{\hat{h}(X_i)}{1 - \hat{h}(X_i)}.$$

But, need  $N > n^2$  to make sure that  $\hat{R} \approx R$ .

# Nonparametric Models: Density Estimation

$$Y_1, \dots, Y_n \sim p(y)$$

$$Y_i \in \mathbb{R}^d$$



# Nonparametric Models: Density Estimation

$$Y_1, \dots, Y_n \sim p(y)$$

$$Y_i \in \mathbb{R}^d$$

We want to estimate the density  $p(y)$  without assuming a model.

## Nonparametric Models: Density Estimation

$$Y_1, \dots, Y_n \sim p(y)$$

$$Y_i \in \mathbb{R}^d$$

We want to estimate the density  $p(y)$  without assuming a model.

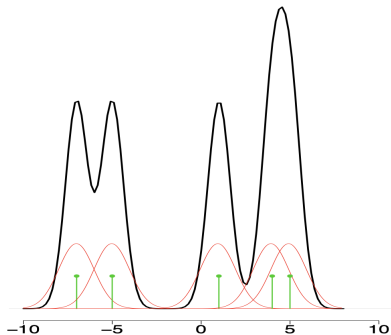
Kernel density estimator:

$$\hat{p}(y) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h^d} K\left(\frac{Y_i - y}{h}\right)$$

where  $K$  is any kernel (example Gaussian) and  $h > 0$  is the bandwidth.

# Kernel Density Estimator

Places a smoothed-out lump of mass of size  $1/n$  over each data point



The shape of the “lumps” is controlled by  $K(\cdot)$ ; their width is controlled by  $h$ .

# Density Estimation

MISE (mean integrated squared error)

$$\mathbb{E}\left[\int (\hat{p}(x) - p(x))^2 dx\right] = \int \text{bias}^2(x) dx + \int \text{var}(x) dx = O(h^4) + O\left(\frac{1}{nh^d}\right)$$

where

$$\text{bias}(x) = \mathbb{E}[\hat{p}(x)] - p(x)$$

# Density Estimation

MISE (mean integrated squared error)

$$\mathbb{E}\left[\int (\hat{p}(x) - p(x))^2 dx\right] = \int \text{bias}^2(x) dx + \int \text{var}(x) dx = O(h^4) + O\left(\frac{1}{nh^d}\right)$$

where

$$\text{bias}(x) = \mathbb{E}[\hat{p}(x)] - p(x)$$

Best  $h = n^{-2/(4+d)}$  gives

$$\text{MISE} = O\left(\left(\frac{1}{n}\right)^{4/(4+d)}\right)$$

# Density Estimation

MISE (mean integrated squared error)

$$\mathbb{E}\left[\int (\hat{p}(x) - p(x))^2 dx\right] = \int \text{bias}^2(x) dx + \int \text{var}(x) dx = O(h^4) + O\left(\frac{1}{nh^d}\right)$$

where

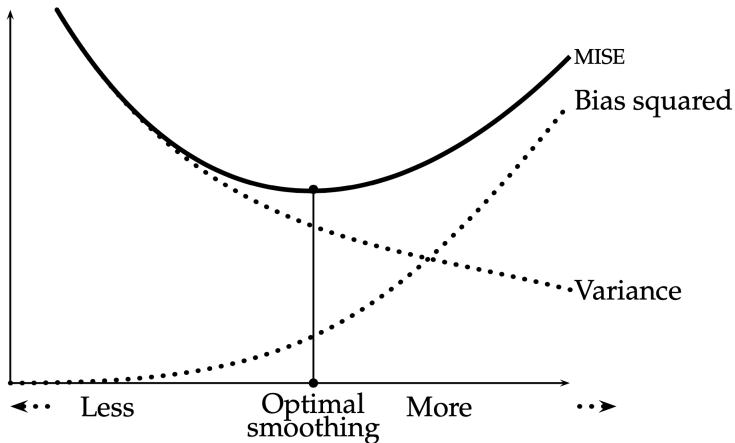
$$\text{bias}(x) = \mathbb{E}[\hat{p}(x)] - p(x)$$

Best  $h = n^{-2/(4+d)}$  gives

$$\text{MISE} = O\left(\left(\frac{1}{n}\right)^{4/(4+d)}\right)$$

Faster rate than histograms.

# Tradeoff



## Choosing $h$ by Cross Validation

Minimize

$$\int (\hat{p}_h(x) - p(x))^2 dx = \int \hat{p}_h^2(x) - 2 \int \hat{p}_h(x)p(x) dx + C$$



## Choosing $h$ by Cross Validation

Minimize

$$\int (\hat{p}_h(x) - p(x))^2 dx = \int \hat{p}_h^2(x) - 2 \int \hat{p}_h(x)p(x) dx + C$$

Estimate by

$$\int \hat{p}_h^2(x) dx - \frac{2}{n} \sum_i \hat{p}_{h,(-i)}(X_i).$$

# Nonparametric Regression (also called machine learning)

$(X_1, Y_1), \dots, (X_n, Y_n)$

If  $Y$  is discrete, we call it classification (or machine learning).

# Nonparametric Regression (also called machine learning)

$$(X_1, Y_1), \dots, (X_n, Y_n)$$

If  $Y$  is discrete, we call it classification (or machine learning).

$$\mu(x) = \mathbb{E}[Y|X = x].$$

# Nonparametric Regression (also called machine learning)

$(X_1, Y_1), \dots, (X_n, Y_n)$

If  $Y$  is discrete, we call it classification (or machine learning).

$\mu(x) = \mathbb{E}[Y|X = x]$ .

Common estimators:

1. kernel
2. local polynomials
3. trees
4. random forests
5. neural nets/deep learning

# Nonparametric Regression

Kernel estimator

$$\hat{\mu}(x) = \frac{\sum_i Y_i K_h(x - X_i)}{\sum_i K_h(x - X_i)} = \sum_i Y_i W_i(x)$$

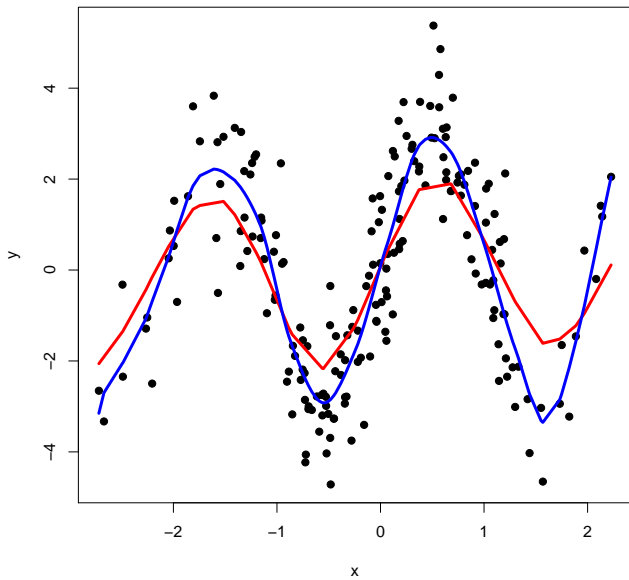
Local polynomial:

$$\mu(u) \approx \beta_0(x)(u - x) + \beta_1(x)(u - x)^2 + \cdots + \beta_p(x)(u - x)^p$$

minimize

$$\sum_i (Y_i - [\beta_0(x)(X_i - x) + \cdots + \beta_p(x)(X_i - x)^p])^2 K_h(x - X_i)$$

# Example



# Trees

Construct an estimator  $\hat{\mu}(x)$  that is a piecewise constant of rectangles:  $R_1, \dots, R_N$  so that

$$\hat{\mu}(x) = \sum_j \bar{Y}_j I(x \in R_j).$$

# Trees

Construct an estimator  $\hat{\mu}(x)$  that is a piecewise constant of rectangles:  $R_1, \dots, R_N$  so that

$$\hat{\mu}(x) = \sum_j \bar{Y}_j I(x \in R_j).$$

The rectangles are constructed recursively:

1. Find best split (which feature? where?)



# Trees

Construct an estimator  $\hat{\mu}(x)$  that is a piecewise constant of rectangles:  $R_1, \dots, R_N$  so that

$$\hat{\mu}(x) = \sum_j \bar{Y}_j I(x \in R_j).$$

The rectangles are constructed recursively:

1. Find best split (which feature? where?)

For example:  $X_7 < 12.2$  versus  $X_7 > 12.2$

# Trees

Construct an estimator  $\hat{\mu}(x)$  that is a piecewise constant of rectangles:  $R_1, \dots, R_N$  so that

$$\hat{\mu}(x) = \sum_j \bar{Y}_j I(x \in R_j).$$

The rectangles are constructed recursively:

1. Find best split (which feature? where?)

For example:  $X_7 < 12.2$  versus  $X_7 > 12.2$

2. Repeat recursively.





# Forests: The Best Off-The-Shelf Method

Subsample data and features:

# Forests: The Best Off-The-Shelf Method

Subsample data and features:

$S_1, \dots, S_N$

## Forests: The Best Off-The-Shelf Method

Subsample data and features:

$S_1, \dots, S_N$

Each subsample  $S_j$  is obtained as follows:

randomly choose a subset of the features (for example,  $d/3$  features):

then draw bootstrap sample of the data:

$(X_1^*, Y_1^*), \dots, (X_n^*, Y_n^*)$

## Forests: The Best Off-The-Shelf Method

Subsample data and features:

$S_1, \dots, S_N$

Each subsample  $S_j$  is obtained as follows:

randomly choose a subset of the features (for example,  $d/3$  features):

then draw bootstrap sample of the data:

$(X_1^*, Y_1^*), \dots, (X_n^*, Y_n^*)$

Fit a tree on each subsample:

$\hat{\mu}_1(x), \dots, \hat{\mu}_N(x)$



## Forests: The Best Off-The-Shelf Method

Subsample data and features:

$$S_1, \dots, S_N$$

Each subsample  $S_j$  is obtained as follows:

randomly choose a subset of the features (for example,  $d/3$  features):

then draw bootstrap sample of the data:

$$(X_1^*, Y_1^*), \dots, (X_n^*, Y_n^*)$$

Fit a tree on each subsample:

$$\hat{\mu}_1(x), \dots, \hat{\mu}_N(x)$$

Then

$$\hat{\mu}(x) = \frac{1}{N} \sum_{j=1}^N \hat{\mu}_j(x)$$

# A Word on Neural Nets/Deep Learning

A nonlinear regression model with many parameters.  
For example:

$$\mu(x) = \alpha + \sum_j \sigma(\beta_j + \gamma_j^T X)$$

# A Word on Neural Nets/Deep Learning

A nonlinear regression model with many parameters.  
For example:

$$\mu(x) = \alpha + \sum_j \sigma(\beta_j + \gamma_j^T X)$$

Requires a lot of tuning and needs a lot of data

## A Word on Neural Nets/Deep Learning

A nonlinear regression model with many parameters.  
For example:

$$\mu(x) = \alpha + \sum_j \sigma(\beta_j + \gamma_j^T X)$$

Requires a lot of tuning and needs a lot of data

Try simpler methods first! In particular, try random forest. Much easier and often works very well.

# Universal Inference (Wasserman, Ramdas and Balakrishnan 2020)

# Universal Inference (Wasserman, Ramdas and Balakrishnan 2020)

Exact inference for hard statistical problems (failure of regularity conditions, nuisance parameters, asymptotics, discrete parameters)

# Universal Inference (Wasserman, Ramdas and Balakrishnan 2020)

Exact inference for hard statistical problems (failure of regularity conditions, nuisance parameters, asymptotics, discrete parameters)

Exact (non asymptotic) coverage. No approximations. No regularity assumptions.

# Background

The usual statistical tests and confidence intervals rely on specific mathematical assumptions.



# Background

The usual statistical tests and confidence intervals rely on specific mathematical assumptions.

Toy example:  $Y_1, \dots, Y_n \sim N(\mu, 1)$

# Background

The usual statistical tests and confidence intervals rely on specific mathematical assumptions.

Toy example:  $Y_1, \dots, Y_n \sim N(\mu, 1)$

Let  $\psi = \mu^2$  and  $\hat{\psi} = \bar{Y}_n^2$ .

# Background

The usual statistical tests and confidence intervals rely on specific mathematical assumptions.

Toy example:  $Y_1, \dots, Y_n \sim N(\mu, 1)$

Let  $\psi = \mu^2$  and  $\hat{\psi} = \bar{Y}_n^2$ .

$\psi = 0$  behaves as  $\chi^2$

$\psi \neq 0$  behaves as Normal

$\psi \approx 0$  behaves as a mixture.

Valid confidence interval?

## Example: Mixtures

$$p(x) = (1 - \lambda)\mathcal{N}(\mu_1, 1) + \lambda\mathcal{N}(\mu_2, 1)$$

(background + signal)

## Example: Mixtures

$$p(x) = (1 - \lambda)N(\mu_1, 1) + \lambda N(\mu_2, 1)$$

(background + signal)

Null (no signal) =

$$\{(\mu_1, \mu_2, \lambda) : \lambda = 0\} = \{(\mu_1, \mu_2, \lambda) : \mu_1 = \mu_2\}$$

## Example: Mixtures

$$p(x) = (1 - \lambda)N(\mu_1, 1) + \lambda N(\mu_2, 1)$$

(background + signal)

Null (no signal) =

$$\{(\mu_1, \mu_2, \lambda) : \lambda = 0\} = \{(\mu_1, \mu_2, \lambda) : \mu_1 = \mu_2\}$$

And when  $\lambda = 0$ , the parameter  $\mu_2$  is not identified.

## Example: Mixtures

$$p(x) = (1 - \lambda)N(\mu_1, 1) + \lambda N(\mu_2, 1)$$

(background + signal)

Null (no signal) =

$$\{(\mu_1, \mu_2, \lambda) : \lambda = 0\} = \{(\mu_1, \mu_2, \lambda) : \mu_1 = \mu_2\}$$

And when  $\lambda = 0$ , the parameter  $\mu_2$  is not identified.

All of our standard machinery fails.

## Other Issues With the Usual Methods

Standard methods are asymptotic (only valid as sample size  $\rightarrow \infty$ )



## Other Issues With the Usual Methods

Standard methods are asymptotic (only valid as sample size  $\rightarrow \infty$ )

Nuisance parameters.

## Other Issues With the Usual Methods

Standard methods are asymptotic (only valid as sample size  $\rightarrow \infty$ )

Nuisance parameters.

Bootstrap: still requires regularity conditions and asymptotics.

## Other Issues With the Usual Methods

Standard methods are asymptotic (only valid as sample size  $\rightarrow \infty$ )

Nuisance parameters.

Bootstrap: still requires regularity conditions and asymptotics.

Bayes: does not provide coverage/error guarantees.

# Universal Inference

Want confidence set  $C_n$  such that

# Universal Inference

Want confidence set  $C_n$  such that

$P(\theta \in C_n) \geq 1 - \alpha$  for all  $\theta$  and all  $n$ . No regularity conditions.  
Works with nuisance parameters.

# Universal Inference

Want confidence set  $C_n$  such that

$P(\theta \in C_n) \geq 1 - \alpha$  for all  $\theta$  and all  $n$ . No regularity conditions.  
Works with nuisance parameters.

Test  $H_0 : \theta \in \Theta_0$ .

Want:  $P(\text{reject})$  under  $H_0$  to be  $\leq \alpha$ . for all  $\theta$  and all  $n$ . No regularity conditions. Works with nuisance parameters.

# Universal Inference

Model  $(p_\theta : \theta \in \Theta)$ .

# Universal Inference

Model ( $p_\theta : \theta \in \Theta$ ).

Split the data into two parts:  $D_0, D_1$ . (We'll get rid the splitting later.)



# Universal Inference

Model ( $p_\theta : \theta \in \Theta$ ).

Split the data into two parts:  $D_0, D_1$ . (We'll get rid the splitting later.)

Get any estimate  $\hat{\theta}$  from  $D_1$ .

# Universal Inference

Model  $(p_\theta : \theta \in \Theta)$ .

Split the data into two parts:  $D_0, D_1$ . (We'll get rid the splitting later.)

Get any estimate  $\hat{\theta}$  from  $D_1$ .

From  $D_0$  construct

$$C = \left\{ \theta : T \geq \alpha \right\}$$

where  $T = \frac{\mathcal{L}(\theta)}{\mathcal{L}(\hat{\theta})}$  and  $L(\theta) = \prod_{i \in D_0} p_\theta(Y_i)$ .

# Universal Inference

Model  $(p_\theta : \theta \in \Theta)$ .

Split the data into two parts:  $D_0, D_1$ . (We'll get rid the splitting later.)

Get any estimate  $\hat{\theta}$  from  $D_1$ .

From  $D_0$  construct

$$C = \left\{ \theta : T \geq \alpha \right\}$$

where  $T = \frac{\mathcal{L}(\theta)}{\mathcal{L}(\hat{\theta})}$  and  $L(\theta) = \prod_{i \in D_0} p_\theta(Y_i)$ .

$C$  is universal:  $P(\theta \in C_n) \geq 1 - \alpha$

for all  $\theta$  and all  $n$ . No asymptotics. No regularity conditions.

## Getting rid of splitting

Split. Get  $T_1$ .

Split. Get  $T_2$ .

Split. Get  $T_3$ .

$\vdots$

Let

$$T = \frac{1}{B} \sum_j T_j$$

and let

$$C = \left\{ \theta : T \geq \alpha \right\}$$

then  $C$  is universal.

## Nuisance Parameters

$$\theta = (\psi, \lambda)$$

$$C = \{ \psi : T \geq \alpha \}$$

where now

$$T = \frac{L(\psi)}{L(\hat{\theta})}$$

where

$$L(\psi) = \sup_{\lambda} L(\psi, \lambda).$$

## Nuisance Parameters

$$\theta = (\psi, \lambda)$$

$$C = \{ \psi : T \geq \alpha \}$$

where now

$$T = \frac{L(\psi)}{L(\hat{\theta})}$$

where

$$L(\psi) = \sup_{\lambda} L(\psi, \lambda).$$

Still valid.

# Hypothesis Testing With Nuisance Parameters

Want to test  $H_0 : \theta \in \Theta_0$ . For example,  $\theta = (\mu, \lambda)$ .

# Hypothesis Testing With Nuisance Parameters

Want to test  $H_0 : \theta \in \Theta_0$ . For example,  $\theta = (\mu, \lambda)$ .

$H_0 : \mu = \mu_0$ .



# Hypothesis Testing With Nuisance Parameters

Want to test  $H_0 : \theta \in \Theta_0$ . For example,  $\theta = (\mu, \lambda)$ .

$H_0 : \mu = \mu_0$ .

$D_1$ : any estimate  $\hat{\theta}$ .

$D_0$ : mle  $\hat{\theta}_0$  under  $H_0$ .

$$T = \frac{L(\hat{\theta})}{L(\hat{\theta}_0)}$$

# Hypothesis Testing With Nuisance Parameters

Want to test  $H_0 : \theta \in \Theta_0$ . For example,  $\theta = (\mu, \lambda)$ .

$H_0 : \mu = \mu_0$ .

$D_1$ : any estimate  $\hat{\theta}$ .

$D_0$ : mle  $\hat{\theta}_0$  under  $H_0$ .

$$T = \frac{L(\hat{\theta})}{L(\hat{\theta}_0)}$$

Reject if  $T \geq 1/\alpha$ .

# Hypothesis Testing With Nuisance Parameters

Want to test  $H_0 : \theta \in \Theta_0$ . For example,  $\theta = (\mu, \lambda)$ .

$H_0 : \mu = \mu_0$ .

$D_1$ : any estimate  $\hat{\theta}$ .

$D_0$ : mle  $\hat{\theta}_0$  under  $H_0$ .

$$T = \frac{L(\hat{\theta})}{L(\hat{\theta}_0)}$$

Reject if  $T \geq 1/\alpha$ .

$P(\text{reject})$  under  $H_0$  is  $\leq \alpha$  for all  $\theta$  and  $n$ , no conditions. Average over splits to get rid of the randomness.

# Hypothesis Testing With Nuisance Parameters

Want to test  $H_0 : \theta \in \Theta_0$ . For example,  $\theta = (\mu, \lambda)$ .

$H_0 : \mu = \mu_0$ .

$D_1$ : any estimate  $\hat{\theta}$ .

$D_0$ : mle  $\hat{\theta}_0$  under  $H_0$ .

$$T = \frac{L(\hat{\theta})}{L(\hat{\theta}_0)}$$

Reject if  $T \geq 1/\alpha$ .

$P(\text{reject})$  under  $H_0$  is  $\leq \alpha$  for all  $\theta$  and  $n$ , no conditions. Average over splits to get rid of the randomness.

$T$  is a p-value

## A Word on Semiparametric Methods

Consider a model like this:

$$p(y; \theta, \gamma)$$

where  $\theta$  is finite dimensional and  $\gamma$  is infinite dimensional.

## A Word on Semiparametric Methods

Consider a model like this:

$$p(y; \theta, \gamma)$$

where  $\theta$  is finite dimensional and  $\gamma$  is infinite dimensional.

We want estimate and confidence interval for  $\theta$ .

## A Word on Semiparametric Methods

Consider a model like this:

$$p(y; \theta, \gamma)$$

where  $\theta$  is finite dimensional and  $\gamma$  is infinite dimensional.

We want estimate and confidence interval for  $\theta$ .

Example:

$$(1 - \lambda)\gamma(x) + \lambda f(x; \theta)$$

where  $\gamma(x)$  is a curve (background).

## A Word on Semiparametric Methods

Consider a model like this:

$$p(y; \theta, \gamma)$$

where  $\theta$  is finite dimensional and  $\gamma$  is infinite dimensional.

We want estimate and confidence interval for  $\theta$ .

Example:

$$(1 - \lambda)\gamma(x) + \lambda f(x; \theta)$$

where  $\gamma(x)$  is a curve (background).

This is called semiparametric inference.



## A Word on Semiparametric Methods

Consider a model like this:

$$p(y; \theta, \gamma)$$

where  $\theta$  is finite dimensional and  $\gamma$  is infinite dimensional.

We want estimate and confidence interval for  $\theta$ .

Example:

$$(1 - \lambda)\gamma(x) + \lambda f(x; \theta)$$

where  $\gamma(x)$  is a curve (background).

This is called semiparametric inference.

Another example:

$$Y = \beta X + f(Z) + \epsilon$$

## A Word on Semiparametric Methods

Consider a model like this:

$$p(y; \theta, \gamma)$$

where  $\theta$  is finite dimensional and  $\gamma$  is infinite dimensional.

We want estimate and confidence interval for  $\theta$ .

Example:

$$(1 - \lambda)\gamma(x) + \lambda f(x; \theta)$$

where  $\gamma(x)$  is a curve (background).

This is called semiparametric inference.

Another example:

$$Y = \beta X + f(Z) + \epsilon$$

This requires a different set of tools.

## A Word on Semiparametric Methods

Consider a model like this:

$$p(y; \theta, \gamma)$$

where  $\theta$  is finite dimensional and  $\gamma$  is infinite dimensional.

We want estimate and confidence interval for  $\theta$ .

Example:

$$(1 - \lambda)\gamma(x) + \lambda f(x; \theta)$$

where  $\gamma(x)$  is a curve (background).

This is called semiparametric inference.

Another example:

$$Y = \beta X + f(Z) + \epsilon$$

This requires a different set of tools.

No time to discuss, but be aware: usual methods don't work!

# Conclusion

Parametric  $\longrightarrow$  Semiparametric  $\longrightarrow$  Nonparametric

# Conclusion

Parametric  $\longrightarrow$  Semiparametric  $\longrightarrow$  Nonparametric

Parametric: MLE, least squares, method of moments, universal

Nonparametric: kernels, trees, forests, deep learning

Semiparametric: other tools

# Conclusion

Parametric  $\longrightarrow$  Semiparametric  $\longrightarrow$  Nonparametric

Parametric: MLE, least squares, method of moments, universal

Nonparametric: kernels, trees, forests, deep learning

Semiparametric: other tools

Estimating likelihoods using classifiers. Great idea but be careful.

# Conclusion

Parametric  $\longrightarrow$  Semiparametric  $\longrightarrow$  Nonparametric

Parametric: MLE, least squares, method of moments, universal

Nonparametric: kernels, trees, forests, deep learning

Semiparametric: other tools

Estimating likelihoods using classifiers. Great idea but be careful.

References:

Wasserman, L. All of Statistics

van der Vaart, A. Asymptotic Statistics.

Freedman, Hastie, Tibshirani. Elements of Statistical Learning.