

Generalization Error Bounds for State Space Models

with an application to economic forecasting

Thesis Proposal

Daniel McDonald

Committee:

Cosma Shalizi, Mark Schervish, Alessandro Rinaldo,
Larry Wasserman, and David N. DeJong

July 22, 2010

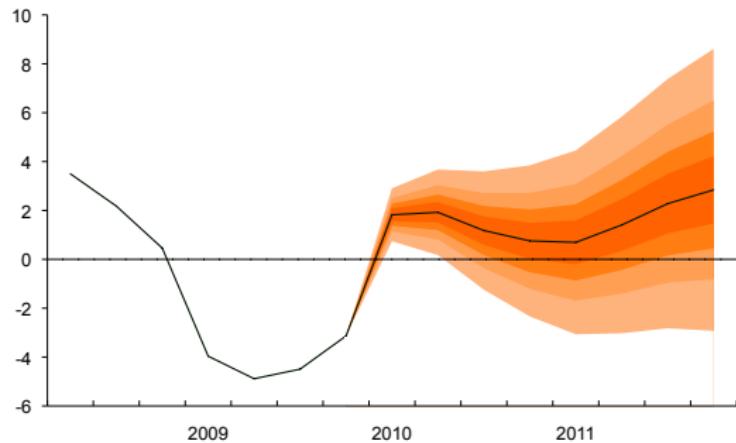
FORECASTING

- Given some data

$$x_1, \dots, x_T \in \mathcal{X}$$

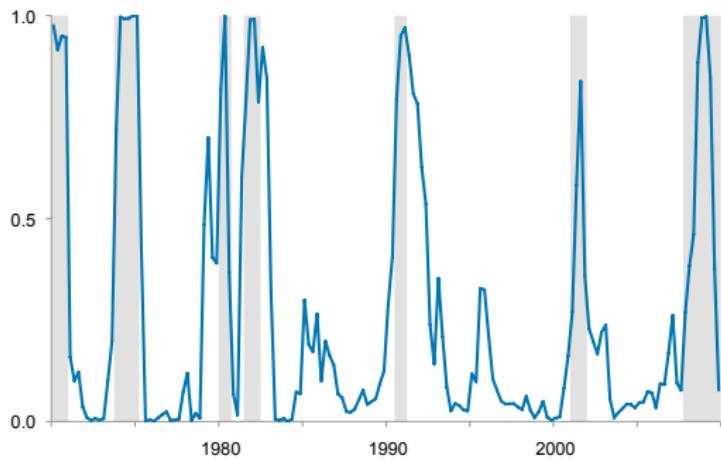
- Want to predict the next data point(s)

$$x_{T+1}, \dots, x_{T+k}$$



Source: Czech National Bank

METHODS OF ECONOMIC FORECASTING

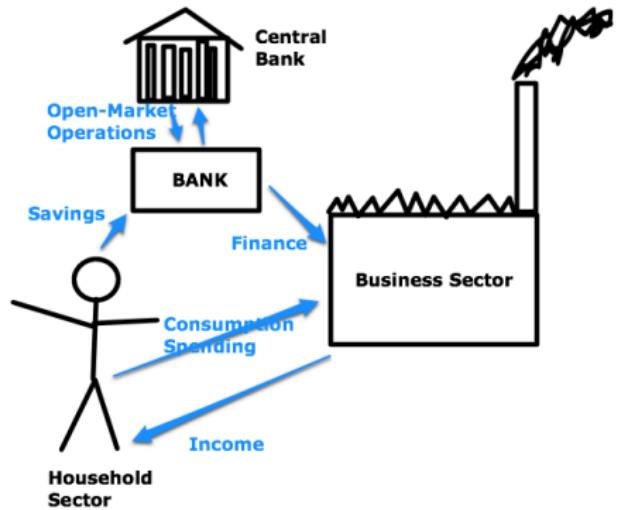


- ARIMA, ARFIMA, GARCH, etc.
- Dynamic Factor Models (Hamilton, Chib, Kim and Nelson, others)
- Systems of Equations models
- Dynamic Stochastic General Equilibrium (DSGE) models

Source: Econbrowser Recession Probabilities

DSGE MODELS

- Most active area of macroeconomic research in the last 30 years
- Arose in response to the Lucas (1976) critique
- Pioneered by Kydland and Prescott (1982)
- Attempt to incorporate “rational behavior” into forecasting models
- Have come under fire for being unable to forecast the financial collapse of 2008–?



Source: Brad DeLong's realization of Daniel Davies' DSGE model

RBC MODEL

- Imagine an infinitely long lived individual who faces the following constrained optimization problem:

$$\max_{c_t, l_t} U = \mathbb{E}_0 \sum_{t=0}^{\infty} \beta^t u(c_t, l_t)$$

$$y_t = z_t q(k_t, n_t)$$

$$1 = n_t + l_t$$

$$y_t = c_t + i_t$$

$$k_{t+1} = i_t + (1 - \delta)k_t$$

$$z_t \sim \text{AR}(1)$$

RELATIONSHIP TO STATE SPACE MODELS

DSGE Model

$$\max_{c_t, l_t} U = \mathbb{E}_0 \sum_{t=0}^{\infty} \beta^t u(c_t, l_t)$$

$$y_t = z_t q(k_t, n_t)$$

$$1 = n_t + l_t$$

$$y_t = c_t + i_t$$

$$k_{t+1} = i_t + (1 - \delta)k_t$$

$$z_t \sim \text{AR}(1)$$

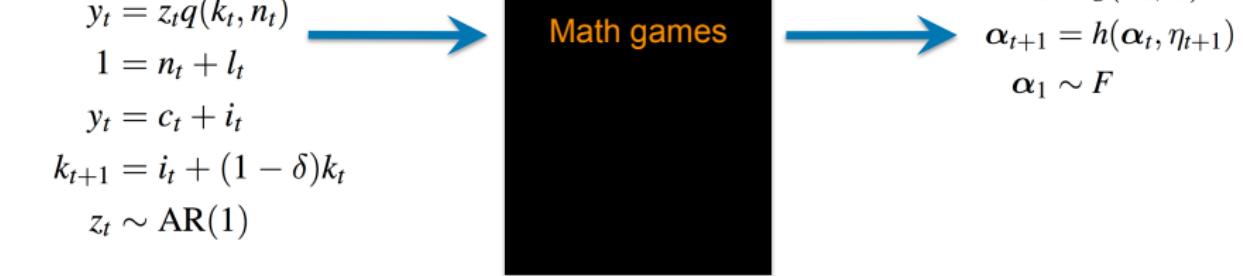
Math games

State Space Model

$$\mathbf{x}_t = g(\boldsymbol{\alpha}_t, \epsilon_t)$$

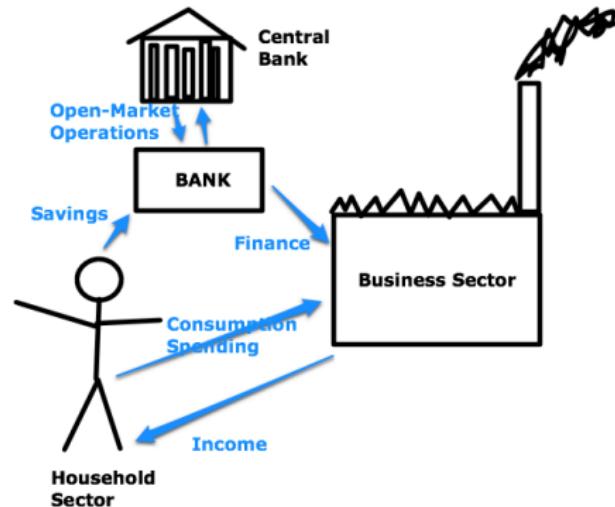
$$\boldsymbol{\alpha}_{t+1} = h(\boldsymbol{\alpha}_t, \eta_{t+1})$$

$$\boldsymbol{\alpha}_1 \sim F$$



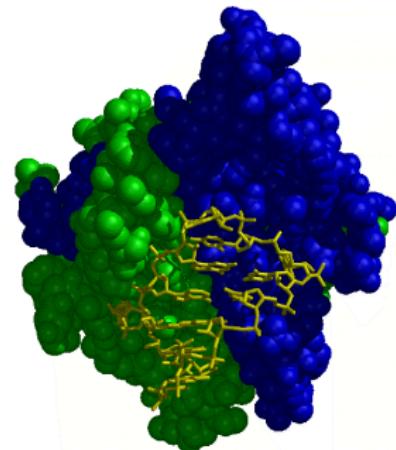
SIMPLICITY/COMPLEXITY

- Unclear if this model is “good”
- Lots of economic arguments
Pro/Con
- What about statistical behavior?
- Overfit/Underfit
- How do predictions compare to other SS models?



STATE SPACE MODELS

- Lots of disciplines use state space models
- Sometimes motivated directly by physical relationships
- How can we measure the forecasting performance?



THESIS PROJECT

PROPOSAL

Develop probabilistic bounds on the prediction error of state space models.

PLAN OF TALK

FORECASTING FRAMEWORK

- 1 Observe training data $D_n = \{(Y_1, X_1), (Y_2, X_2), \dots, (Y_n, X_n)\}$ from some stochastic process μ
- 2 Choose model class \mathcal{F} from which to construct predictors, e.g. AR(p), DSGE, regression, wavelets, etc.
- 3 Use a loss function $\ell(Y, f(X))$ to measure performance of candidate predictors $f \in \mathcal{F}$
- 4 Estimate the model using D_n , to produce \hat{f} , your proposed forecasting model

GENERALIZATION ERROR

- Want to control the generalization error, or risk, of chosen predictor \hat{f}

$$R(\hat{f}) = \mathbb{E}_\mu[\ell(Y_0, \hat{f}(X_0)) \mid D_n]$$

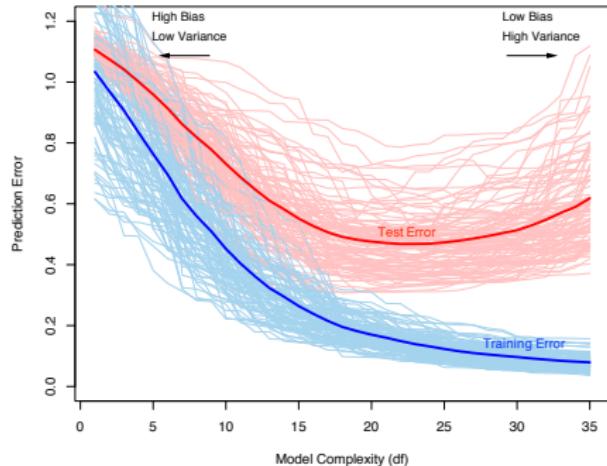
- But the stochastic process μ is unknown
- Usually estimate $R(\hat{f})$ with training error

$$R_n(\hat{f}) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i, \hat{f}(X_i))$$

- This is a poor estimate of $R(\hat{f})$

TRAINING ERROR

- Often $R_n(\hat{f}) < R(\hat{f})$
- Complicated classes \mathcal{F} lead to small $R_n(\hat{f})$
- Simple classes \mathcal{F} lead to large $R_n(\hat{f})$
- Model comparisons using $R_n(\hat{f})$ lead to choosing overly complex \mathcal{F} —overfitting
- Related to the bias/variance tradeoff



Source: Hastie, Tibshirani, and Friedman *The Elements of Statistical Learning*

RISK

$$R(\hat{f}) = \mathbb{E}_{\mu}[\ell(Y_0, \hat{f}(X_0)) \mid D_n]$$

- Estimation of $R(\hat{f})$ is a hard problem since μ is unknown
- Instead, derive probabilistic upper bounds
- These bounds depend on \mathcal{F} —one needs to characterize the complexity of different function classes, as well as the algorithms used to choose \hat{f}
- Many complexity measures—VC dimension, covering numbers, algorithmic stability, and Rademacher complexity

RADEMACHER COMPLEXITY

DEFINITION

Define the Rademacher complexity of a function class \mathcal{F} as

$$\mathfrak{R}(\mathcal{F}) = \mathbb{E}_X \mathbb{E}_\sigma \left[\sup_{f \in \mathcal{F}} \left| \frac{2}{n} \sum_{i=1}^n \sigma_i f(x_i) \right| \right],$$

where σ_i are iid and $\mathbb{P}(\sigma_i = 1) = \mathbb{P}(\sigma_i = -1) = \frac{1}{2}$.

- Measures the maximum covariance between the predictions and random noise—how closely can some $f \in \mathcal{F}$ fit garbage?
- Gives tight bounds
- Removing \mathbb{E}_X gives empirical Rademacher complexity
- Somewhat ugly, hard to calculate

IID RESULTS FOR DEPENDENT DATA

- IID bounds require n independent training samples
- Applying IID results to dependent data requires knowing the number of “independent” samples ($< n$)

INTUITION

For dependent data, knowledge of a few data points gives information about other data points. The marginal benefit of observing new data is therefore less than if the data were independent

- Measures of **mixing** “adjust” the IID bounds

MIXING

DEFINITION

Let $\{X_i\}_{i=1}^{\infty}$ be random variables generated from a stationary process. Let $\mathbb{P}_{t \otimes t+m}$ denote the joint distribution between $\{X_i\}_{i=1}^t$ and $\{X_i\}_{i=t+m}^{\infty}$ with \mathbb{P}_t and \mathbb{P}_{t+m} the associated marginals. Then the β -mixing coefficient is given by

$$\beta(m) = \sup_t \|\mathbb{P}_t \otimes \mathbb{P}_{t+m} - \mathbb{P}_{t \otimes t+m}\|_{TV}.$$

- Sufficient condition¹ to apply IID results to dependent data is

$$\beta(m) \xrightarrow{m \rightarrow \infty} 0$$

- If this holds, the process is β -mixing

¹ Yu 1994

TIME SERIES BOUNDS

THEOREM

Let \mathcal{H} be the space of losses bounded above by M . Then given a sample from a stationary β -mixing distribution, for all $m, a > 0$ with $2ma = n$ and $\eta > 2(a - 1)\beta(m)$, then for all $f \in \mathcal{F}$, with probability at least $1 - \eta$,

$$R(f) < R_a(f) + \mathfrak{R}_a(\mathcal{H}) + M\sqrt{\frac{\log 2/\eta'}{2a}}$$

with $\eta' = \eta - 2(a - 1)\beta(m)$.

Source: Mohri and Rostamizadeh 2009

IMPLICATIONS

THEOREM

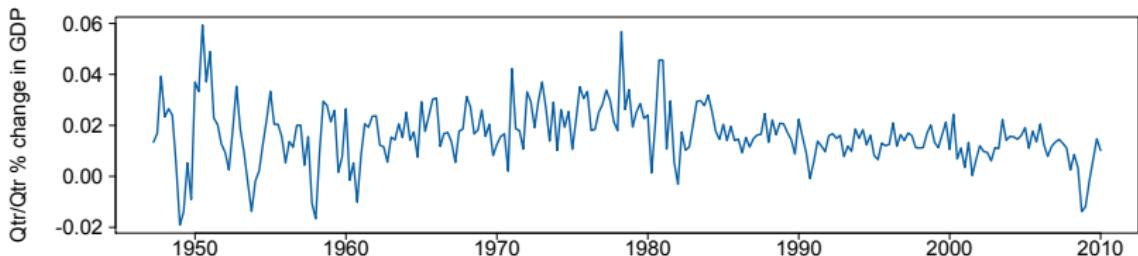
$$R(f) < R_a(f) + \mathfrak{R}_a(\mathcal{H}) + M \sqrt{\frac{\log 2/\eta'}{2a}}$$

- The effective sample size is not n but a
- The empirical risk is based on a data points separated by a distance $2m$
- Faster decay in $\beta(m)$ means more ‘independent’ samples, smaller third term
- Second term is Rademacher complexity of the loss space
- Can substitute empirical Rademacher complexity with slight modifications
- M is an upper bound for the loss

THIS BOUND DOES NOT WORK!

- Do not yet know Rademacher complexity for **general** state space models
- The β -mixing rate is assumed known—the bound depend of properties of unknown data generating process
- The bounded loss requirement is unpleasant for regression problems

EXAMPLE



- 1 GDP data from 1947–2010 ($n = 252$)
- 2 Fit an AR(2) model
- 3 Assuming it is true, calculate β mixing rate
- 4 Calculate empirical Rademacher complexity (uses properties of regularized kernel regression)
- 5 Use squared error loss, truncated at M

With probability at least 0.85,

$$R(\hat{f}) \leq 1 \times 10^{-4} + 0.07\sqrt{M} + 1.03M$$

WHAT HAPPENED?

- For all M , $M < \text{bound}$
 - ➡ Bound holds, but it is trivial
- 0.85 confidence level is strange
 - ➡ Due to $\beta(m)$
- There is not much data
 - ➡ ~ 14 ‘independent’ data points!
- Second and third terms have the wrong order of magnitude

COMPONENTS OF PROJECT

PROPOSAL

Develop generalization error bounds for state space models to control the predictive risk and allow for model selection.

- 1 Estimate mixing behavior from data, adjust to accommodate extra uncertainty
- 2 Measure complexity of state space models
- 3 Modify bounds to allow for unbounded loss
- 4 Investigate resampling methods

ESTIMATION OF MIXING RATES

Three avenues for estimation so far:

1 ‘Plug-in’ method

- Estimate the two densities using standard techniques
- Numerical integration to get the total variation distance
- **Badly biased (upward)**

2 ‘Information’ method

- Can bound β -mixing: $\beta(m) \leq \sqrt{I(m)}$
- $I(m)$ is just mutual information
- **Bound may be loose, some methods have very bad behavior for nearly independent distributions**

3 ‘Copula’ method

- Total variation distance is invariant under 1-1 transformations
- Problem reduces to estimating the copula and a simpler numerical integration

When solving a given problem, try to avoid solving a more general problem as an intermediate step. —Vladimir Vapnik

COMPLEXITY CHARACTERIZATION

1 Rademacher complexity

- Gives tight bounds
- Can calculate for stationary AR (regularized kernel problem)
- Not so obvious for more complicated models (VAR, linear SS, nonlinear SS)

2 VC-dimension

- More likely to work for complicated models—Can likely calculate simply
- Available bounds not in terms of VC-dim

UNBOUNDED LOSS

- Requires altering the existing bounds
- Need different concentration inequalities
- These require **different** assumptions on the data generating process
 - Result from Jiang (2010)—Need to control

$$\int_c^{\infty} |X_t| dp(x)$$

- Result from Vapnik—Need to control

$$\sup_{f \in \mathcal{F}} \frac{\left(\int \ell^p(Y_0, f(X_0)) d\mu \right)^{1/p}}{\int \ell(Y_0, f(X_0)) d\mu}$$

ALTERNATIVE METHODS

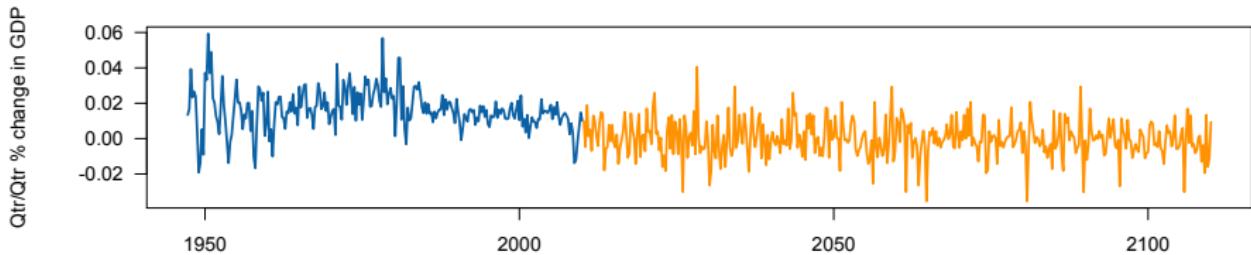
1 Cross-validation/Subset prediction

- Most common ‘validation’ method in economic forecasting literature
- Some theory
- Models are complex/not much data
- Models are slow to estimate
- Modelers have observed held out data thoroughly

2 Bootstrapping

- Need theory for time series generalization error
- Models are slow to estimate
- Could work...

BOOTSTRAP ERROR BOUND



- 1 Fit an AR(2) model and calculate the training error
- 2 Repeat B times:
 - Bootstrap a new series which is several times longer than the original
 - Fit a model to initial portion and calculate the in-sample risk
 - Calculate the prediction risk on the remainder
 - Store the difference between the in-sample and prediction risks
- 3 Find the 95% and 85% of the distribution of over-fits. Add this to the training error

$$R(\hat{f}) \leq 1.5 \times 10^{-4}$$

$$R(\hat{f}) \leq 1.35 \times 10^{-4}$$

CONCLUSIONS

PROPOSAL

Develop generalization error bounds for state space models to control the predictive risk and allow for model selection.

- 1 Estimate mixing behavior from data, adjust the bound(s) to accommodate this additional uncertainty
- 2 Measure complexity of state space models
- 3 Modify bounds to allow for unbounded loss, refine bounded loss results (bounded bounds!)
- 4 Investigate resampling methods

THE END

Thanks for coming

STATIONARY AR

- Bound the Rademacher complexity of the class of models

$$\mathcal{F}_p = \left\{ \varphi_1, \dots, \varphi_p : x_t = \sum_{i=1}^p \varphi_i x_{t-i} \text{ and } x_t \text{ is stationary} \right\}$$

- Stationarity requires the roots of $p(z) = z^p + \varphi_1 z^{p-1} + \dots + \varphi_p$ lie inside the complex unit disc.
- Can show that a sufficient condition is¹

$$\|\varphi\|_2^2 \leq \sum_{i=1}^p \binom{p}{i}^2 = \binom{2p}{p} - 1$$

- Bounds the norm

¹ Fam and Meditch 1978

STATIONARY AR

- Ordinary linear regressions can be written as kernel regressions. Let

$$\alpha_i = (\mathbf{X}(\mathbf{X}'\mathbf{X})^{-2}\mathbf{X}'\mathbf{Y})_i$$
$$k(\mathbf{X}_i, \mathbf{X}_j) = \mathbf{X}_i \mathbf{X}'_j,$$

where \mathbf{X} is the $n \times p$ design matrix, \mathbf{Y} are the responses, and \mathbf{X}_i is the i^{th} row of the design matrix.

- Requiring $\sum_{i,j} \alpha_i \alpha_j k(\mathbf{X}_i, \mathbf{X}_j) \leq \gamma^2$
- Corresponds $||\hat{\beta}^{OLS}||_2^2 \leq \gamma^2$, or ridge regression

STATIONARY AR

$$\mathcal{F}_p \subseteq \overline{\mathcal{F}_p} = \left\{ \varphi_1, \dots, \varphi_p : x_t = \sum_{i=1}^p \varphi_i x_{t-i} \text{ and } \|\varphi\|_2^2 \leq \binom{2p}{p} - 1 \right\}$$

Allows application of kernel regularized result¹

$$\mathfrak{R}(\mathcal{F}_p) \leq \mathfrak{R}(\overline{\mathcal{F}_p}) \leq \frac{2}{\sqrt{n}} \sqrt{\left(\binom{2p}{p} - 1 \right) \mathbb{E} \mathbf{X}_1 \mathbf{X}_1'}$$

$$\mathfrak{R}_n(\mathcal{F}_p) \leq \mathfrak{R}_n(\overline{\mathcal{F}_p}) \leq \frac{2}{\sqrt{n}} \sqrt{\left(\binom{2p}{p} - 1 \right) \frac{1}{n} \sum_{t=i}^n \mathbf{X}_i \mathbf{X}_i'}$$

¹ Bartlett and Mendelson 2002

EXAMPLE

THEOREM

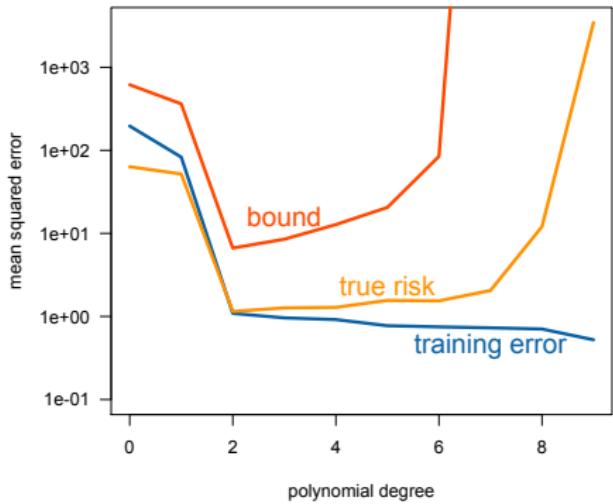
Let $h = d_{\mathcal{F}}/n$ where $d_{\mathcal{F}}$ is the VC dimension of \mathcal{F} . Then with probability at least $1 - \eta$

$$R(\hat{f}) \leq R_n(\hat{f}) \times \left(1 - \sqrt{h - h \ln h - \frac{\ln \eta}{n}}\right)_+^{-1}$$

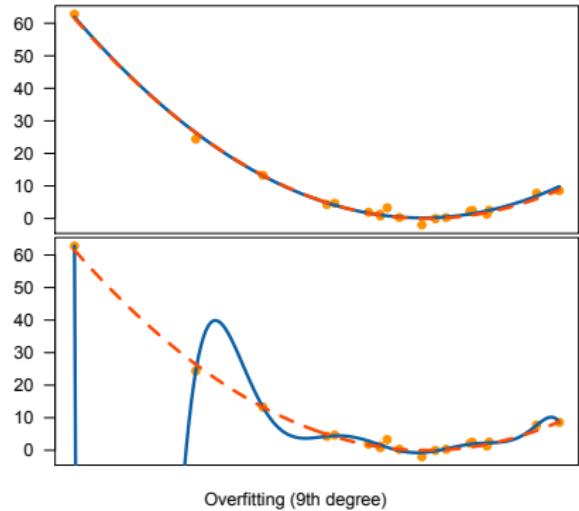
- Applies to all nonnegative loss functions ℓ and distributions μ satisfying certain conditions¹
- Is easy to apply for many classes \mathcal{F}
- Can use for model selection

¹ Vapnik

EXAMPLE



Optimal model (2nd degree)



Overfitting (9th degree)

ON-LINE LEARNING RESULTS

Cesa-Bianchi, Conconi, and Gentile 2004

- Nice results about generalization error (rather than regret) in on-line framework
- **Main idea**—Take the hypothesis generated at the end of each time step, calculate the loss on future data, add complexity penalty. Hypothesis with smallest penalized empirical risk is used. Its risk is bounded easily.
- Suffers from a few issues
 - 1 Applies to IID data
 - 2 Bounded loss
 - 3 Bootstrap/CV flavor is computationally difficult

We close by mentioning an important open question: Is it possible to extend the results...to the case when the examples are not independent; e.g., when they are generated by a stationary process?

Could be a nice problem to tackle along the way

ON-LINE LEARNING RESULTS

Rakhlin, Sridharan, and Tewari 2010

- **Main idea**—Provides conditions which ensure online learnability
- The online learning problem is cast as a game played on rooted binary trees
- Defines complexity measures for bounding the value of the game (regret) of an online learner
- **Important result for time series**—Finite fat-shattering dimension implies universal uniform convergence (an analogue of Glivenko-Cantelli for dependent data)
- Bounded function classes with finite fat shattering dimension are online learnable
- In this case, various complexity measures (including the authors' sequential Rademacher complexity) are within $O(\log^{3/2} n)$

ESTIMATING β -MIXING

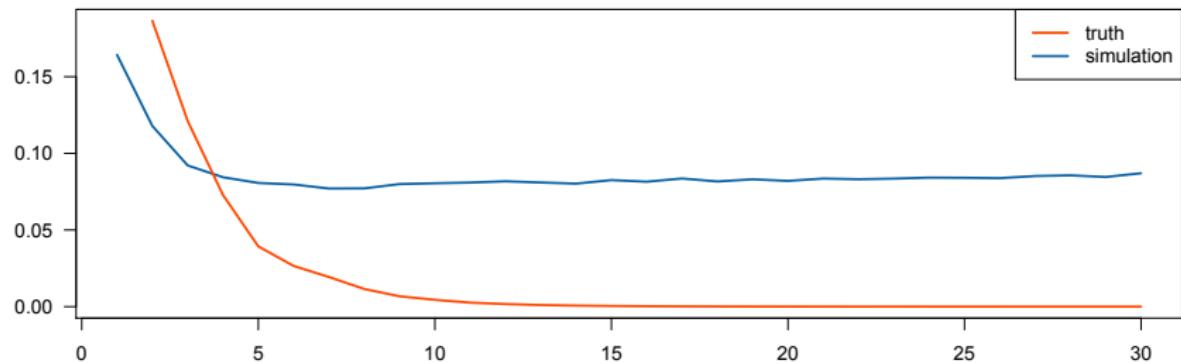
We want to estimate

$$\begin{aligned}\beta(m) &= \sup_t ||\mathbb{P}_t \otimes \mathbb{P}_{t+m} - \mathbb{P}_{t \otimes t+m}||_{TV} \\ &= \frac{1}{2} \int \int |p_t(x) \otimes p_{t+m}(y) - p_{t \otimes t+m}(x, y)| dx dy\end{aligned}$$

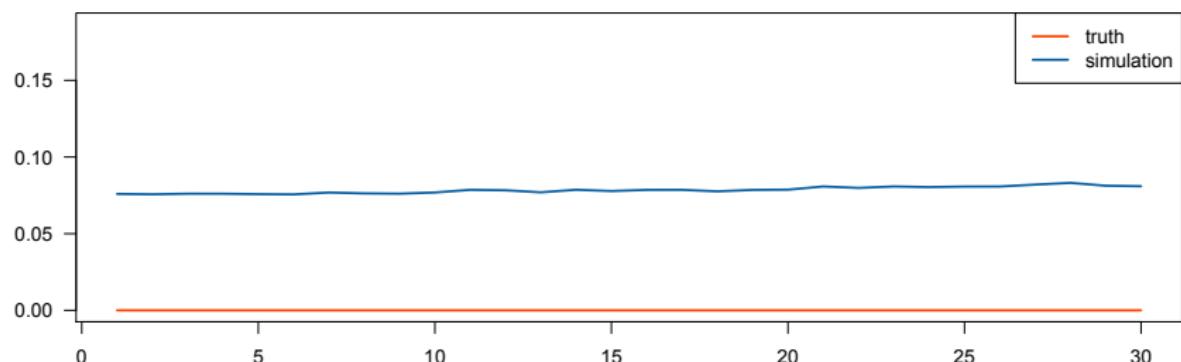
- Since we have assumed stationarity $p_t(x) = p_{t+m}(y)$
- Thus there is only a univariate density and a bivariate density to estimate

ESTIMATING β -MIXING—RESULTS

Plug-in simulation AR(2)



Plug-in simulation N(0,1)



ESTIMATING INFORMATION MIXING

We want to estimate

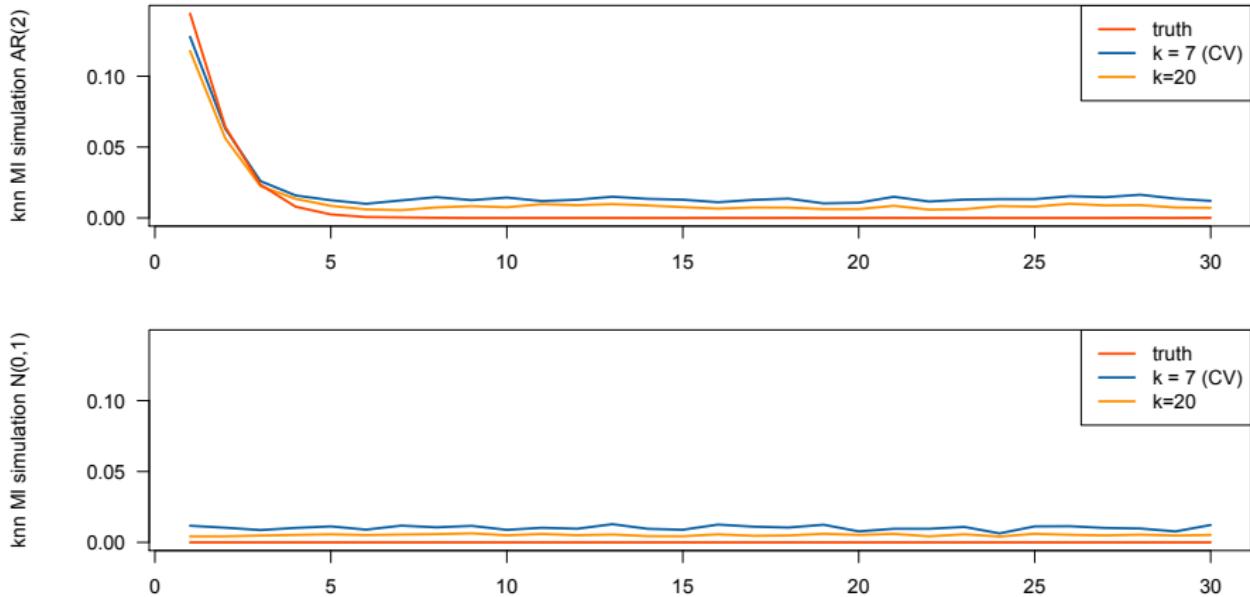
$$I(m) = \int \log \frac{p_{t \otimes t+m}(x, y)}{p_t(x)p_{t+m}(y)} dp_{t \otimes t+m}(x, y)$$

- Could use plug in estimator
- Tried k -nn type estimate of Kraskov, Stögbauer, and Grassberger 2004

$$\widehat{I(X, Y)} = \psi(k) - \frac{1}{n} \sum_{i=1}^n [\psi(n_x(i) + 1) + \psi(n_y(i) + 1)] + \psi(n)$$

where $\psi(\cdot)$ is the digamma function

ESTIMATING INFORMATION MIXING—RESULTS



INSIDE THE DSGE BLACK BOX (LINEAR)

- 1 Find the first order conditions for the system
- 2 Holding the stochastic process constant, derive the steady state of the system
- 3 Perform a (first order) Taylor series expansion around the steady state
- 4 Use the method of [Sims 2001](#) to convert to linear state space form
(basically a QZ decomposition)

INSIDE THE DSGE BLACK BOX (NON-LINEAR)

- 1 Find the first order conditions for the system
- 2 Holding the stochastic process constant, derive the steady state of the system
- 3 Use either of
 - Perturbation methods (higher order Taylor series expansions) using Schmitt-Grohé and Uribe 2004
 - Expand the offending nonlinear functions in some basis (projection methods) using Judd 1998

BLAH BLAH

$$\mathbf{x}_t = g(\boldsymbol{\alpha}_t, \epsilon_t)$$

$$\boldsymbol{\alpha}_{t+1} = h(\boldsymbol{\alpha}_t, \eta_{t+1})$$

$$\boldsymbol{\alpha}_1 \sim F$$