

Model Selection and Computational Oracle Inequalities for Large Scale Problems

Peter Bartlett

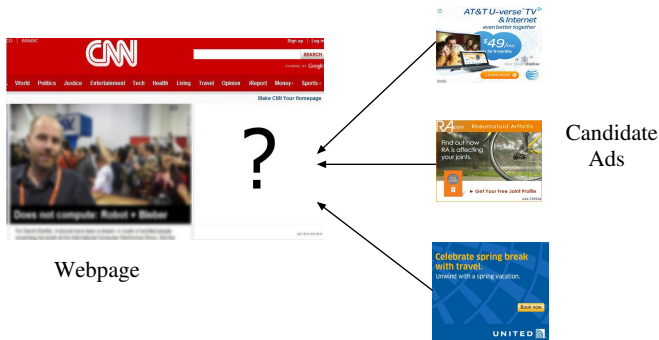
University of California at Berkeley

Queensland University of Technology

MMDS, July 2012

Model Selection: Click Prediction for Online Ads

- Predict **click/no-click** given advertisement and webpage.
- Training data from past click logs.
- e.g.: Logistic regression with 16M parameters.



- 1 Model selection
 - Prediction problems
 - Model selection
 - Oracle inequalities
 - Fast rates
 - Oracle inequalities with fast rates
- 2 Large scale model selection
 - Computational oracle inequalities
 - Fast rates
- 3 Summary and open problems

Prediction in a Probabilistic Setting

- i.i.d. Z, Z_1, \dots, Z_n from \mathcal{Z} .

Example:

$Z = (X, Y)$, where

X contains features of advertisement and webpage.

Y indicates 'click' or 'no click.'

Prediction in a Probabilistic Setting

- i.i.d. Z, Z_1, \dots, Z_n from \mathcal{Z} .
- Use data Z_1, \dots, Z_n to choose $f_n \in \mathcal{A}$ with small **risk**,

$$L(f_n) = \mathbf{E}\ell(f_n, Z).$$

Here, $\ell : \mathcal{A} \times \mathcal{Z} \rightarrow \mathbb{R}^+$ is a **loss** function.

Examples:

- 1 $\ell(f, (x, y)) = (y - f(x))^2$ for $y \in \{\pm 1\}$ and $f : \mathcal{X} \rightarrow [0, 1]$: estimate probability of click.
- 2 $\ell(f, (x, y)) = 1[y \neq f(x)]$ for $y \in \{\pm 1\}$ and $f : \mathcal{X} \rightarrow \{\pm 1\}$: decide most likely class label.
- 3 $\ell(f, z) = -\log f(z)$: estimate probability density.

Prediction in a Probabilistic Setting

- i.i.d. Z, Z_1, \dots, Z_n from \mathcal{Z} .
- Use data Z_1, \dots, Z_n to choose $f_n \in \mathcal{A}$ with small **risk**,

$$L(f_n) = \mathbf{E} \ell(f_n, Z).$$

Here, $\ell : \mathcal{A} \times \mathcal{Z} \rightarrow \mathbb{R}^+$ is a **loss** function.

- Choose f_n from a class F .

Examples:

- 1 Linear regression:

$$F = \{x \mapsto \theta'x : \theta \in \mathbb{R}^p\}.$$

- 2 Logistic regression:

$$F = \{x \mapsto \sigma(\theta'x) : \theta \in \mathbb{R}^p\}.$$

- 3 Support vector machines:

$$F = \text{Reproducing kernel Hilbert space.}$$

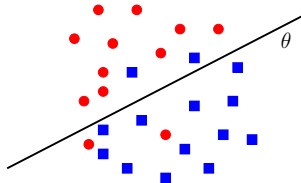
Approximation-Estimation Trade-Off

- Define the *Bayes risk*, $L^* = \inf_f L(f)$, where the infimum is over measurable f .
- We can decompose the excess risk as

$$L(\hat{f}) - L^* = \underbrace{\left(L(\hat{f}) - \inf_{f \in F} L(f) \right)}_{\text{estimation error}} + \underbrace{\left(\inf_{f \in F} L(f) - L^* \right)}_{\text{approximation error}}.$$

- **Model selection**: automatically choose F to optimize this trade-off.

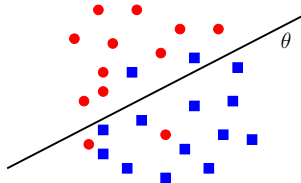
Example 1: Norm of a linear predictor



- Many linear classification algorithms minimize:

$$\min_{\theta \in \mathbb{R}^p} \sum_{i=1}^n \ell(y_i, \langle \theta, x_i \rangle) \quad \text{subject to} \quad \|\theta\|_2 \leq r.$$

Example 1: Norm of a linear predictor



- Many linear classification algorithms minimize:

$$\min_{\theta \in \mathbb{R}^p} \sum_{i=1}^n \ell(y_i, \langle \theta, x_i \rangle) \quad \text{subject to} \quad \|\theta\|_2 \leq r.$$

- Estimation and approximation errors depend on the bound r
- Often select from a grid $r_1 \leq r_2 \leq r_3 \leq \dots$

Example 2: Feature selection

- $\theta \in \mathbb{R}^p$, select subset of $\{1, 2, \dots, p\}$ where $\theta_i \neq 0$

Example 2: Feature selection

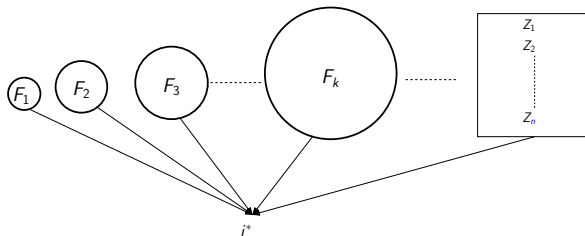
- $\theta \in \mathbb{R}^p$, select subset of $\{1, 2, \dots, p\}$ where $\theta_i \neq 0$
- Examples:
 - Natural language: Unigrams \prec Bigrams $\prec \dots \prec n$ -grams
 - Function fitting: polynomial degree, Fourier basis dim, ...
 - Computer vision: hierarchy of wavelet filters

Example 2: Feature selection

- $\theta \in \mathbb{R}^p$, select subset of $\{1, 2, \dots, p\}$ where $\theta_i \neq 0$
- Examples:
 - Natural language: Unigrams \prec Bigrams $\prec \dots \prec n$ -grams
 - Function fitting: polynomial degree, Fourier basis dim, ...
 - Computer vision: hierarchy of wavelet filters
- Approximation and estimation errors depend on dimensionality.

The Model Selection Problem

- Hierarchy of model classes, F_1, F_2, F_3, \dots
- Data Z_1, Z_2, \dots, Z_n



Want i^* that optimizes estimation-approximation trade-off

$$L(\hat{f}_i) - L(f^*) = \underbrace{(L(\hat{f}_i) - \inf_{f \in F_i} L(f))}_{\text{Estimation error}} + \underbrace{(\inf_{f \in F_i} L(f) - L(f^*))}_{\text{Approximation error}}$$

The Model Selection Problem

Given function classes F_1, F_2, \dots , use the data Z_1, \dots, Z_n to choose $\hat{f} \in \bigcup_i F_i$ that gives a good trade-off between the approximation error and the estimation error.

Example: *Complexity-penalized model selection.*

$$f_n^i = \arg \min_{f \in F_i} L_n(f),$$

$$\hat{f} = \text{minimizer of } L_n(f_n^i) + \gamma_i(n),$$

where $\gamma_i(n)$ is a *complexity penalty* and L_n is the empirical risk:

$$L_n(f) = \frac{1}{n} \sum_{i=1}^n \ell(f, Z_i).$$

A Simple Oracle Inequality

Theorem

Suppose that we have risk bounds for each F_i : w.p. $1 - \delta$,

$$\sup_{f \in F_i} |L(f) - L_n(f)| \leq \gamma_i(n) + c \sqrt{\frac{\log 1/\delta}{n}}.$$

If \hat{f} is chosen via complexity regularization:

$$f_n^i = \arg \min_{f \in F_i} L_n(f), \quad \hat{f} = \text{minimizer of } L_n(f_n^i) + \gamma_i(n),$$

then with probability $1 - \delta$,

$$L(\hat{f}) \leq \min_i \left(\inf_{f \in F_i} L(f) + 2\gamma_i(n) + c \sqrt{\frac{\log 1/\delta + \log i}{n}} \right).$$

A Simple Oracle Inequality

- Notice that, for each F_i satisfying

$$\sup_{f \in F_i} |L(f) - L_n(f)| \leq \gamma_i(n) + c \sqrt{\frac{\log 1/\delta}{n}},$$

we have
$$L(f_n^i) \leq \inf_{f \in F_i} L(f) + 2\gamma_i(n) + c \sqrt{\frac{\log 1/\delta}{n}}.$$

- But complexity regularization gives \hat{f} satisfying

$$L(\hat{f}) \leq \min_i \left(\inf_{f \in F_i} L(f) + 2\gamma_i(n) + c \sqrt{\frac{\log 1/\delta + \log i}{n}} \right).$$

- Thus, \hat{f} gives a near-optimal trade-off between the approximation error and the (bound on) estimation error, with only a $\log i$ penalty.

Fast Rates

These oracle inequalities rely on uniform convergence: bounds on

$$\sup_{f \in F_i} |L(f) - L_n(f)|.$$

Typical fluctuations are of the order

$$|L(f) - L_n(f)| = O\left(\frac{1}{\sqrt{n}}\right).$$

In some cases, these rates cannot be improved, and additive penalties that scale as

$$\sup_{f \in F_i} |L(f) - L_n(f)| = \Omega\left(\frac{1}{\sqrt{n}}\right)$$

give optimal oracle inequalities.

Fast Rates

However, in many cases, we can obtain faster rates.
e.g., with high probability, for all $f \in F$,

$$L(f) - L(f^*) \leq 2(L_n(f) - L_n(f^*)) + O\left(\frac{\log n}{n}\right),$$

where $L(f^*) = \min_{f \in F} L(f)$. In these cases, choosing

$$\hat{f} = \arg \min_{f \in F} L_n(f)$$

gives $L(\hat{f}) \leq L(f^*) + O(\log n/n)$.

Examples: Convex losses [Lee, B., Williamson, 1998; B., Jordan, McAuliffe, 2006], classification with low noise [Mammen and Tsybakov, 2004; Tsybakov, 2004].

Oracle Inequalities with Fast Rates for Complexity Regularization

It turns out that we can use complexity regularization to exploit these faster rates, provided the F_i are ordered by inclusion.

Theorem (B., 2008)

For $F_1 \subseteq F_2 \subseteq \dots$ and $\gamma_1(n) \leq \gamma_2(n) \leq \dots$, if

$$\sup_i \sup_{f \in F_i} (L(f) - L(f_i^*) - 2(L_n(f) - L_n(f_i^*)) - \gamma_i(n)) \leq 0,$$

$$\sup_i \sup_{f \in F_i} (L_n(f) - L_n(f_i^*) - 2(L(f) - L(f_i^*)) - \gamma_i(n)) \leq 0,$$

$$\text{then } L(\hat{f}) \leq \inf_i (L(f_i^*) + 9\gamma_i(n)),$$

where \hat{f} minimizes $L_n(f_i^*) + 7\gamma_i(n)/2$ and $f_i^* = \arg \min_{f \in F_i} L(f)$.

Oracle Inequalities with Fast Rates for Complexity Regularization

This is *striking*:

- $L_n(f_n^i)$ fluctuates on a scale $1/\sqrt{n}$.
- But adding a tiny penalty $\gamma_i(n) = O(\log n/n)$ gives $L(\hat{f})$ within $O(\log n/n)$ of the best!

The explanation: the fluctuations for different F_i are correlated, because the empirical minimizers are chosen using the *same data* and the F_i are ordered by inclusion.

- 1 Model selection
 - Prediction problems
 - Model selection
 - Oracle inequalities
 - Fast rates
 - Oracle inequalities with fast rates
- 2 Large scale model selection
 - Computational oracle inequalities
 - Fast rates
- 3 Summary and open problems

Large Scale Data Analysis

Joint work with Alekh Agarwal, John Duchi and Clément Levrard.

Observation:

For many prediction problems, the amount of data available is *effectively unlimited*.

Large Scale Data Analysis

Observation:

For many prediction problems, the amount of data available is *effectively unlimited*.

Natural language processing:

Spelling correction

Google Linguistics Data

Consortium n -gram corpus:

10^{11} sentences.

muamar gadafi|

About 10,200 results (0.25 seconds)

Did you mean: [muammar gaddafi](#)

Large Scale Data Analysis

Observation:

For many prediction problems, the amount of data available is *effectively unlimited*.

Computer vision: Captions

Facebook:

10^{11} photos.



United Nations Secretary General **Kofi Annan** stands with U.N. Security Council President and U.S. Ambassador to the U.N. **John D. Negroponte** as Annan ...



North Korean leader **Kim Jong Il**, and Russian President **Vladimir Putin** walk after talks in Vladivostok, Friday, Aug. 23, 2002. North Korean leader Kim ...

Large Scale Data Analysis

Observation:

For many prediction problems, the amount of data available is *effectively unlimited*.

- Information retrieval: Web search
- Natural language processing: Spelling correction
- Computer vision: Captions

Large Scale Data Analysis

Observation:

For many prediction problems, performance is limited by *computational resources*, not sample size.

- Information retrieval: Web search
- Natural language processing: Spelling correction
- Computer vision: Captions

Large Scale Data Analysis

Example:

- Peter Norvig, “Internet-Scale Data Analysis”:
On a spelling correction problem, trivial prediction rules, estimated with a massive dataset perform much better than complex prediction rules (which allow only a dataset of modest size).
- Given a limited computational budget,
what is the best trade-off?
That is, should we spend our computation on gathering more data, or on estimating richer prediction rules?

- 1 Model selection
 - Prediction problems
 - Model selection
 - Oracle inequalities
 - Fast rates
 - Oracle inequalities with fast rates
- 2 Large scale model selection
 - Computational oracle inequalities
 - Fast rates
- 3 Summary and open problems

Computation versus sample size

Recall: *Complexity-penalized model selection.*

$$f_n^i = \arg \min_{f \in F_i} L_n(f),$$

$$\hat{f} = \text{minimizer of } L_n(f_n^i) + \gamma_i(n),$$

where $\gamma_i(n)$ is a *complexity penalty* and L_n is the empirical risk:

$$L_n(f) = \frac{1}{n} \sum_{i=1}^n \ell(f, Z_i).$$

Computation versus sample size

- Complexity regularization involves computation of the empirical risk minimizer for each F_i :

$$f_n^i = \arg \min_{f \in F_i} L_n(f), \quad \hat{f} = \text{minimizer of } L_n(f_n^i) + \gamma_i(n),$$

So computation typically grows *linearly with number of classes*.

- The **oracle inequality** gives the best trade-off *for a given sample size*:

$$L(\hat{f}) \leq \min_i \left(\inf_{f \in F_i} L(f) + 2\gamma_i(n) + c \sqrt{\frac{\log 1/\delta + \log i}{n}} \right).$$

Scaling of penalties with computation

Recall

$\gamma_i(n)$ is the complexity penalty for the class F_i with sample size n .

Scaling of penalties with computation

Recall

$\gamma_i(n)$ is the complexity penalty for the class F_i with sample size n .

Define

$p_i(T)$ as the complexity penalty for the class F_i with computational budget T .

computation $T \implies$ sample size $n_i(T)$ for F_i

We set $p_i(T) = \gamma_i(n_i(T))$.

Scaling of penalties with computation

Define

$p_i(T)$ as the complexity penalty for the class F_i with computational budget T .

In more detail:

with computation T , we can ensure that, with high probability,

$$\sup_{f \in F_i} |L(f) - L_{n_i(T)}(f)| \leq \gamma_i(n_i(T)),$$

hence

$$L(f_{n_i(T)}^i) \leq \inf_{f \in F_i} L(f) + O(p_i(T)).$$

Scaling of penalties with computation

Define

$p_i(T)$ as the complexity penalty for the class F_i with computational budget T .

Our goal: A computational oracle inequality:
 \hat{f} compares favorably with each model, estimated using the entire computational budget.

$$L(\hat{f}) \leq \min_i \left(\underbrace{\inf_{f \in F_i} L(f) + O(p_i(T))}_{\text{c.f. estimate } f \text{ using the entire budget}} \right).$$

Scaling of penalties with computation

Define

$p_i(T)$ as the complexity penalty for the class F_i with computational budget T .

Our goal: A computational oracle inequality:
 \hat{f} compares favorably with each model, estimated using the entire computational budget.

$$L(\hat{f}) \leq \min_i \left(\underbrace{\inf_{f \in F_i} L(f)}_{\text{c.f. estimate } f \text{ using almost the entire budget}} + O\left(p_i\left(\frac{T}{\log T}\right)\right) \right).$$

Naïve solution: grid search

- Allocate budget T/K to each model F_1, \dots, F_K .
- Use a sample of size $n_i(T/K)$ for F_i .
- Choose

$$f_{n_i}^i = \arg \min_{f \in F_i} L_{n_i}(f),$$

$$\hat{f} = \text{minimizer of } L_{n_i}(f_{n_i}^i) + \gamma_i(n_i).$$

- Satisfies oracle inequality

$$L(\hat{f}) \leq \min_i \left(\inf_{f \in F_i} L(f) + p_i \left(\frac{T}{K} \right) \right).$$

Model selection from nested classes

- Suppose that the models are ordered by inclusion:

$$F_1 \subseteq F_2 \subseteq \dots$$

- Examples:

- $F_i = \{f_\theta : \theta \in \mathbb{R}^p, \|\theta\| \leq r_i\}, r_1 \leq r_2 \leq \dots.$
- $F_i = \{f_\theta : \theta \in \mathbb{R}^{p_i}, \|\theta\| \leq 1\}, p_1 \leq p_2 \leq \dots.$

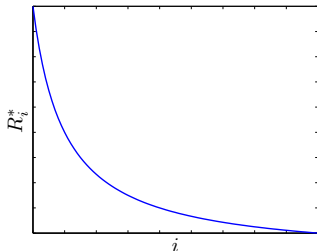
- Suppose that we have risk bounds for each F_i : w.p. $1 - \delta$,

$$\sup_{f \in F_i} |L(f) - L_n(f)| \leq \gamma_i(n) + c \sqrt{\frac{\log 1/\delta}{n}}.$$

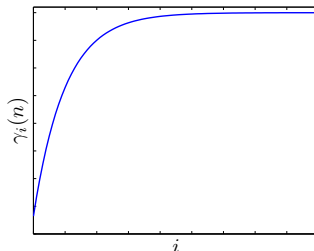
Exploiting structure of nested classes

Want to exploit monotonicity of risks and penalties

Excess risk, $R_i^* = \inf_{f \in F_i} L(f) - L^*$:

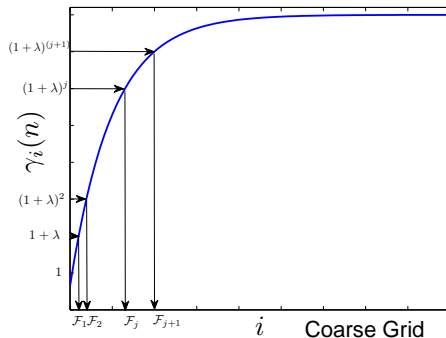


Penalty, $\gamma_i(n)$:



Coarse grid sets

- Want to spend computation on only few classes.
- Use monotonicity to interpolate for the rest.
- Partition based on penalty values.



Coarse grids for model selection

Assume

- 1 Loss is bounded:

$$\ell(f, Z) \in [0, B].$$

- 2 Computation grows at least linearly with sample size:

$$n_1(T) = O(T).$$

- 3 Penalty decreases no faster than $1/n$:

$$\gamma_1(n) = \Omega\left(\frac{1}{n}\right).$$

Coarse grids for model selection

Then

- We can ignore F_i with $\gamma_i(n_i(T)) > B$.
- We can cover all smaller classes with a **coarse grid** of size $s = O(\log(BT))$.

Definition (Coarse grid)

For $S \subseteq \mathbb{N}$, a set $\hat{S} \subseteq S$ is a **coarse grid** of size s for S if $|\hat{S}| = s$ and for each $i \in S$ there is an index $j \in \hat{S}$ such that

$$\gamma_i \left(n_i \left(\frac{T}{s} \right) \right) \leq \gamma_j \left(n_i \left(\frac{T}{s} \right) \right) \leq 2\gamma_i \left(n_i \left(\frac{T}{s} \right) \right).$$

Coarse grids for model selection

Then

- We can ignore F_i with $\gamma_i(n_i(T)) > B$.
- We can cover all smaller classes with a **coarse grid** of size $s = O(\log(BT))$.
- Include a new class only after penalty increases sufficiently.
- $s = \log \left(\frac{B}{\gamma_1(n_1(T))} \right) = O(\log BT)$ suffices.

Complexity regularization on a coarse grid

Given a coarse grid \hat{S} with cardinality s :

- ① Allocate budget T/s to each class in S .
- ② Choose

$$f^i = \arg \min_{f \in F_i} L_{n_i(T/s)}(f)$$

$$\hat{f} = \arg \min_{f \in \{f^j : j \in \hat{S}\}} L_{n_j(T/s)}(f) + \gamma_j \left(n_j \left(\frac{T}{s} \right) \right).$$

Complexity regularization on a coarse grid

Theorem

For a nested hierarchy satisfying the uniform convergence bounds, with high probability,

$$\begin{aligned} L(\hat{f}) &\leq \min_i \left\{ \inf_{f \in F_i} L(f) + O \left(\gamma_i \left(n_i \left(\frac{T}{s} \right) \right) \right) \right\} \\ &\leq \min_i \left\{ \inf_{f \in F_i} L(f) + O \left(p_i \left(\frac{T}{\log T} \right) \right) \right\} \end{aligned}$$

- *Computational cost of model selection scales logarithmically with T .*

- 1 Model selection
 - Prediction problems
 - Model selection
 - Oracle inequalities
 - Fast rates
 - Oracle inequalities with fast rates
- 2 Large scale model selection
 - Computational oracle inequalities
 - Fast rates
- 3 Summary and open problems

Computational Oracle Inequalities?

Can we obtain computational oracle inequalities with fast rates?

Computational Oracle Inequalities?

Can we obtain computational oracle inequalities with fast rates?

Previous Algorithm

Given a coarse grid \hat{S} with cardinality s :

- 1 Allocate budget T/s to each class in S .
- 2 Choose

$$f^i = \arg \min_{f \in F_i} L_{n_i(T/s)}(f)$$

$$\hat{f} = \arg \min_{f \in \{f^j : j \in \hat{S}\}} L_{n_j(T/s)}(f) + \gamma_j \left(n_j \left(\frac{T}{s} \right) \right).$$

Computational Oracle Inequalities?

Previous Algorithm

Given a coarse grid \hat{S} with cardinality s :

- 1 Allocate budget T/s to each class in S .
- 2 Choose

$$f^i = \arg \min_{f \in F_i} L_{n_i(T/s)}(f)$$

$$\hat{f} = \arg \min_{f \in \{f_j : j \in \hat{S}\}} L_{n_j(T/s)}(f) + \gamma_j \left(n_j \left(\frac{T}{s} \right) \right).$$

Obstacle: The oracle inequality relies on the use of the *same data*. But to best use our computational budget, we should gather *more* data for simpler classes.

Algorithm for Fast Rates

Given a coarse grid \hat{S} with cardinality s :

- 1 Allocate budget T/s to each class in S .
- 2 Choose

$$f^i = \arg \min_{f \in F_i} L_{n_i}(T/s^2)(f)$$

- 3 Define \hat{f} as the f^i with the largest index i such that for all smaller j ,

$$L_{n_i}(f^i) + \gamma_i(n_i) \leq \inf_{f \in F_j} L_{n_j}(f) + \gamma_j(n_j).$$

The *same data* is used in comparing f^i with functions from smaller classes.

Computational Oracle Inequalities

Theorem

For a nested hierarchy exhibiting fast rates, with high probability,

$$L(\hat{f}) \leq \min_i \left\{ \inf_{f \in F_i} L(f) + O \left(p_i \left(\frac{T}{\log^2 T} \right) \right) \right\}.$$

- 1 Model selection
 - Prediction problems
 - Model selection
 - Oracle inequalities
 - Fast rates
 - Oracle inequalities with fast rates
- 2 Large scale model selection
 - Computational oracle inequalities
 - Fast rates
- 3 Summary and open problems

Open problems

- For nested hierarchies, the analysis relied on a coarse multiplicative cover of the penalty values. If the penalties are data-dependent, when is this approach possible?
- What other structures on function classes lead to good computational oracle inequalities?
- How do computational constraints affect the optimal performance in other estimation problems?

Summary

- For large-scale problems, data is cheap but computation is precious.
- Computational oracle inequalities for model selection: select a near-optimal model without wasting much computation on other models.
- A *nested* complexity hierarchy ensures cost logarithmic in computational budget.
- Faster rates are sometimes possible: More complicated complexity regularization schemes ensure cost polylogarithmic in computational budget.
- If not nested, cost of model selection is linear in size of hierarchy.