# Distribution-Free Prediction

## Perspectives from a recovering Bayesian

Dave Zachariah

Talk @ ERNSI 2017

dave.zachariah@it.uu.se

# Background

📄 Wågberg, J., Zachariah, D., Schön, T.B., Stoica, P.: Prediction performance after learning in Gaussian process regression, AISTATS, 2017.

📄 Zachariah, D., Stoica, P., and Schön, T.B.: Online Learning for Distribution-Free Prediction, submitted, 2017.

📄 Zachariah, D. and Stoica, P.: Cramér-Rao Bounds in Statistical Inference, in preparation.

# Aims

1. Subjectivist-Bayesian vs. Frequentist interpretations
2. Bayesian paradigm is powerful but sometimes misleading
3. Distribution-free approach to tackle challenges
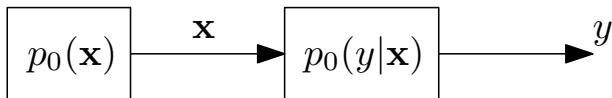
# Prediction problem

## Data-generating process

$$\boxed{p_0(\mathbf{x})} \xrightarrow{\ \mathbf{x}\ } \boxed{p_0(y|\mathbf{x})} \longrightarrow y$$

Figure : Process with inputs $\mathbf{x}$ and outputs $y$

$$\text{Dataset } \mathcal{D} = \big\{(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_n, y_n)\big\}$$

# Data-generating process



$$p_0(\mathbf{x}) \xrightarrow{\ \mathbf{x}\ } p_0(y|\mathbf{x}) \longrightarrow y$$

Figure : Process with inputs $\mathbf{x}$ and outputs $y$

Example #1:

- $\mathbf{x}$ spatial coordinates
- $y$ ozone density
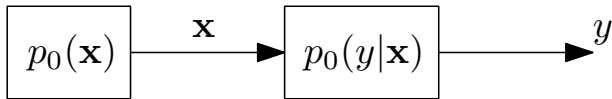- $n = 17\ 340$ samples
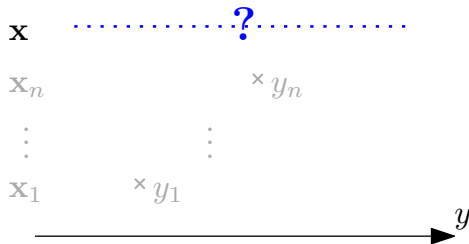
# Data-generating process



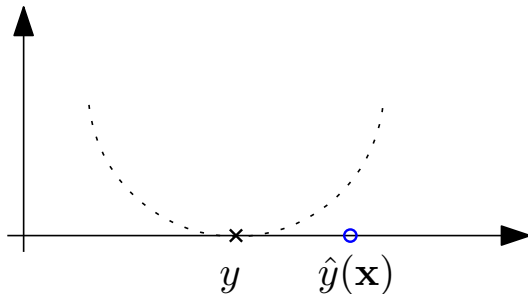Figure : Process with inputs $\mathbf{x}$ and outputs $y$

Example #2:

- $\mathbf{x}$ individual background covariates
- $y$ weekly wage
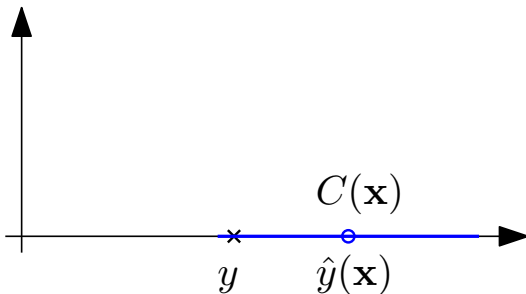- $n = 219\ 673$ samples

# Prediction method as a decision rule



- ▶ New sample $(\mathbf{x}, y) \sim p_0(\mathbf{x}, y)$
- ▶ Given $\mathbf{x}$ predict $y$ using $\mathcal{D}$

dave.zachariah@it.uu.se

**Prediction method as a decision rule**



Decision-rule $\widehat{y}(\mathbf{x})$ has risk $\mathcal{R} = \mathrm{E}\Big[|y - \widehat{y}(\mathbf{x})|^2\Big]$
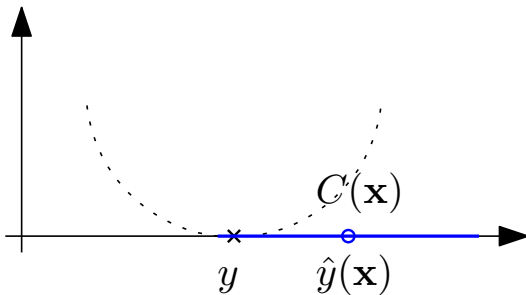
# Uncertainty region



Uncertainty region

$$C(\mathbf{x}) = \left\{ y' \ : \ d\left(y', \widehat{y}(\mathbf{x})\right) \leq \kappa \right\}$$
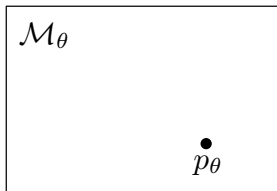
constructed using $\mathcal{D}$

 dave.zachariah@it.uu.se

**Decisions under uncertainty**



To proceed we specify a model $p_\theta$ of unknown $p_0$

# Specifying a model class

$\mathcal{M}_\theta$

$\overset{\bullet}{p_\theta}$
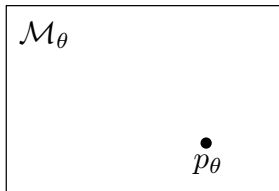
Data-generating process:

$$p_0(\mathbf{x}, y, \mathcal{D})$$

Model:

$$p_\theta(\mathbf{x}, y, \mathcal{D})$$

# Model class

$$\boxed{\begin{array}{l} \mathcal{M}_\theta \\[2.5em] \qquad\qquad \overset{\bullet}{\phantom{.}}_{\displaystyle p_\theta} \end{array}}$$
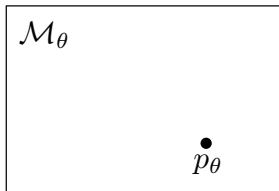
Data-generating process:

$$p_0(\mathbf{x}, y, \mathcal{D})$$

Model (w/ marginalized latent variables):

$$p_\theta(\mathbf{x}, y, \mathcal{D}) = \int p_\theta(\mathbf{x}, y, \mathbf{z}, \mathcal{D}) \, d\mathbf{z}$$

# Model class

$$
\boxed{
\begin{array}{l}
\mathcal{M}_\theta \\[3em]
\qquad\qquad\quad \bullet \\
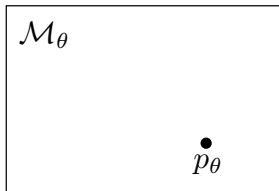\qquad\qquad\quad p_\theta
\end{array}
}
$$

Data-generating process:

$$p_0(\mathbf{x}, y, \mathcal{D})$$

Model (factorized form):

$$p_\theta(\mathbf{x}, y, \mathcal{D}) = p_\theta(\mathbf{x}, \mathcal{D})\, p_\theta(y|\mathbf{x}, \mathcal{D})$$
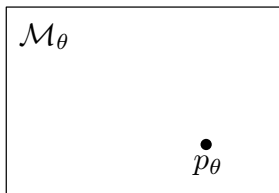
# Model class

$$\boxed{\begin{array}{l} \mathcal{M}_\theta \\[2.5cm] \qquad\qquad \overset{\textstyle\bullet}{p_\theta} \end{array}}$$

Example: i.i.d. samples

$$\mathcal{M}_\theta = \left\{ \, p_\theta(\mathbf{x}, \mathcal{D}) \, p_\theta(y|\mathbf{x}, \mathcal{D}) \, : \, \mu_\theta(\mathbf{x}) = \boldsymbol{\phi}^\top(\mathbf{x})\boldsymbol{\theta} \, \right\}$$

with feature vector $\boldsymbol{\phi}(\mathbf{x})$

## Model class



$\mathcal{M}_\theta$

$\overset{\bullet}{p_\theta}$

Example: Gaussian process

$$\mathcal{M}_\theta = \left\{ p_\theta(\mathbf{x}, \mathcal{D}) \, p_\theta(y|\mathbf{x}, \mathcal{D}) \text{ Gaussian } : \, \mu_\theta(\mathbf{x}), k_\theta(\mathbf{x}, \mathbf{x}') \right\}$$
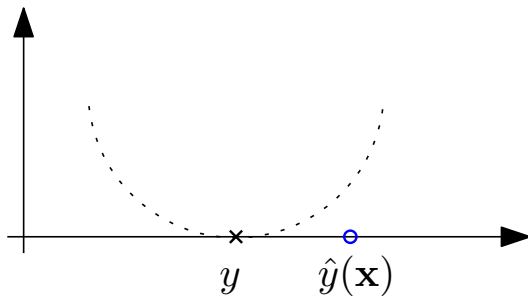
with posterior mean and covariance functions using $\mathcal{D}$

# Ideal case

$$\boxed{\begin{array}{l} \mathcal{M}_\theta \\[2em] \qquad\qquad \overset{\bullet}{p_\theta = p_0} \end{array}}$$

## Optimal decision-rule



### Bound on risk

$$\mathcal{R} \geq \mathrm{E}_{x,\mathcal{D}}\Big[ \mathrm{Var}[y|\mathbf{x}, \mathcal{D}] \Big]$$

Attained by decision rule $\widehat{y}(\mathbf{x}) = \mu_{\theta}(\mathbf{x})$

## Uncertainty region


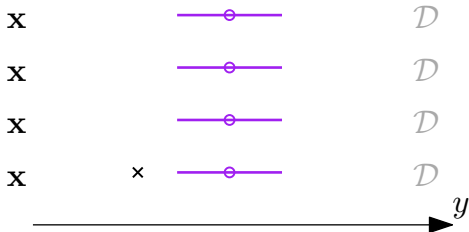
Since $\sigma_\theta^2(\mathbf{x}) = \mathrm{Var}[y|\mathbf{x}, \mathcal{D}]$, construct uncertainty region

$$C_\theta(\mathbf{x}) = \left\{ \, y' \, : |y' - \mu_\theta(\mathbf{x})| \leq \kappa \sigma_\theta(\mathbf{x}) \, \right\}$$

# Uncertainty region: Credibility

## Credibility

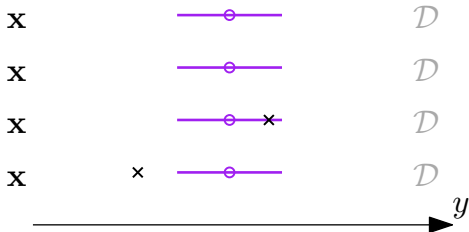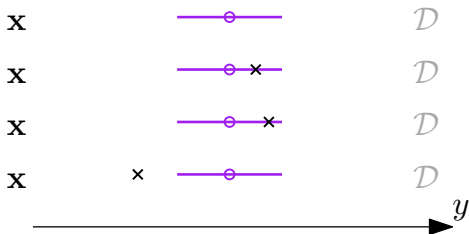$$\Pr\{y \in C_\theta(\mathbf{x}) \mid \mathbf{x}, \mathcal{D}\} \geq 1 - \kappa^{-2}$$



Individualized decision-making and beliefs

# Uncertainty region: Credibility

## Credibility

$$\Pr\{y \in C_\theta(\mathbf{x}) \mid \mathbf{x}, \mathcal{D}\} \ \geq \ 1 - \kappa^{-2}$$
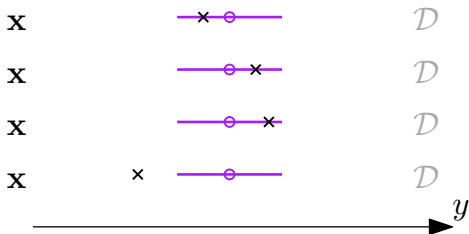


Individualized decision-making and beliefs

# Uncertainty region: Credibility

## Credibility

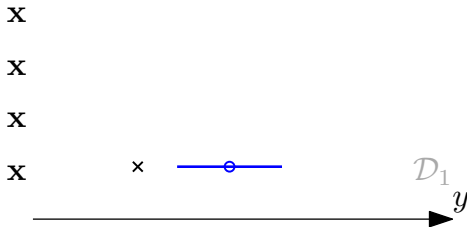$$\Pr\{y \in C_\theta(\mathbf{x}) \mid \mathbf{x}, \mathcal{D}\} \geq 1 - \kappa^{-2}$$



Individualized decision-making and beliefs

dave.zachariah@it.uu.se

# Uncertainty region: Credibility

## Credibility

$$\Pr\big\{y \in C_\theta(\mathbf{x}) \mid \mathbf{x}, \mathcal{D}\big\} \ \geq \ 1 - \kappa^{-2}$$



Individualized decision-making and beliefs

## Uncertainty region: Confidence

### Confidence

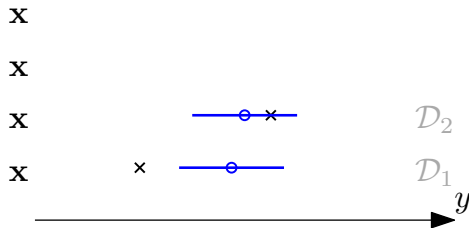$$\Pr\{y \in C_\theta(\mathbf{x}) \mid \mathbf{x}\} \geq 1 - \kappa^{-2}$$

$\mathbf{x}$

$\mathbf{x}$

$\mathbf{x}$

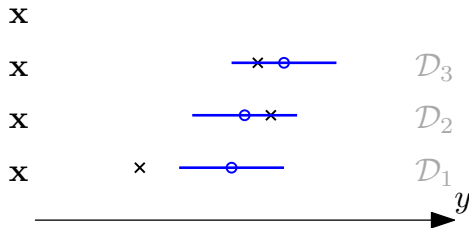$\mathbf{x}$     ×    ———o———     $\mathcal{D}_1$
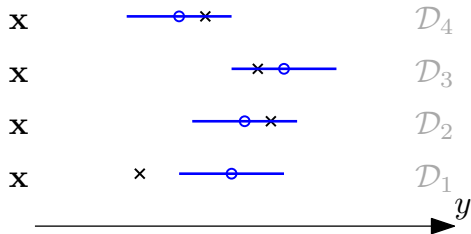
$y$

Reproducible decision-making and coverage

## Uncertainty region: Confidence

### Confidence

$$\Pr\{y \in C_\theta(\mathbf{x}) \mid \mathbf{x}\} \ \geq \ 1 - \kappa^{-2}$$



Reproducible decision-making and coverage

# Uncertainty region: Confidence

### Confidence

$$\Pr\big\{y \in C_\theta(\mathbf{x}) \mid \mathbf{x}\big\} \geq 1 - \kappa^{-2}$$



Reproducible decision-making and coverage
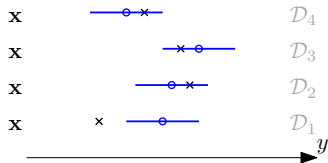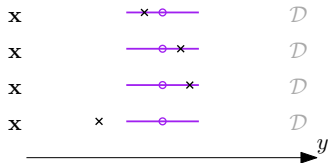
# **Uncertainty region: Confidence**

### Confidence

$$\Pr\{y \in C_\theta(\mathbf{x}) \mid \mathbf{x}\} \geq 1 - \kappa^{-2}$$



Reproducible decision-making and coverage
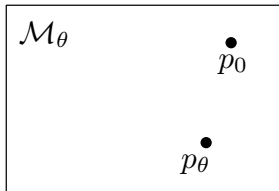
# Subjectivist vs. Frequentist interpretations



|  | Subjectivist: | Frequentist: |
|--|--|--|
|  | Belief | Coverage |
|  | Credibility | Confidence |
|  | Individualized | Reproducible |

# Well-specified case

# Defining the best model



Figure : Divergence $\Delta$ of $p_\theta$ from $p_0$

## Defining the best model



Figure : Divergence $\Delta$ of $p_\theta$ from $p_0$

Best model:

$$\boldsymbol{\theta}_\star = \arg\min_{\boldsymbol{\theta}} \ \Delta(\boldsymbol{\theta}) \qquad \text{e.g.} \quad \Delta(\boldsymbol{\theta}) = \widehat{\mathsf{E}}\Big[|y_i - \mu_{\boldsymbol{\theta}}(\mathbf{x}_i)|^2\Big]$$

## Defining the best model



Figure : Divergence $\Delta$ of $p_\theta$ from $p_0$

Best model:

$$\boldsymbol{\theta}_\star = \arg\min_{\boldsymbol{\theta}} \; \Delta(\boldsymbol{\theta}) \qquad \text{e.g.} \quad \Delta(\boldsymbol{\theta}) = \mathrm{E}_{y|X}\left[\ln \frac{p_0(\mathbf{y}|\mathbf{X})}{p_\theta(\mathbf{y}|\mathbf{X})}\right]$$
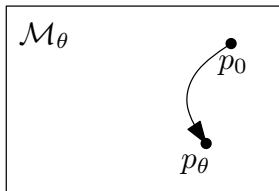
# Defining the best model



Figure : Divergence $\Delta$ of $p_\theta$ from $p_0$

Well-specified case: $\quad p_0 \in \mathcal{M}_\theta \quad \Rightarrow \quad p_{\theta_\star} = p_0$

# Defining the best model



Figure : Divergence $\Delta$ of $p_\theta$ from $p_0$

Learned model $\widehat{\boldsymbol{\theta}}$ that approaches $\boldsymbol{\theta}_\star$ as $n$ grows

# Bound on risk when best model is unknown



$$y \qquad \hat{y}(\mathbf{x})$$

### Bound in Gaussian case

If bias is invariant w.r.t. $\boldsymbol{\theta}_\star$:

$$\mathcal{R} \geq \mathrm{E}_{x,X}\Big[ \mathrm{Var}[y|\mathbf{x}, \mathbf{X}] + \mathbf{g}^\top \mathbf{J}^{-1}\mathbf{g} \Big],$$
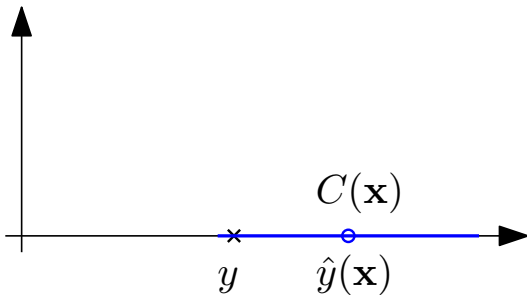
where $\mathbf{J}$ is a Fisher information matrix.

## Uncertainty region: unknown properties



Constructed as before but evaluted at $\widehat{\boldsymbol{\theta}}$:

$$C_{\widehat{\theta}}(\mathbf{x}) = \left\{ y' : |y' - \mu_{\widehat{\theta}}(\mathbf{x})| < \kappa \sigma_{\widehat{\theta}}(\mathbf{x}) \right\}$$

dave.zachariah@it.uu.se

# Uncertainty region: unknown properties



Credibility:

$$\Pr\left\{y \in C_{\widehat{\theta}}(\mathbf{x}) \mid \mathbf{x}, \mathcal{D}\right\} \quad ?$$

Confidence:

$$\Pr\left\{y \in C_{\widehat{\theta}}(\mathbf{x}) \mid \mathbf{x}\right\} \quad ?$$

Approximations possible when $n \to \infty$
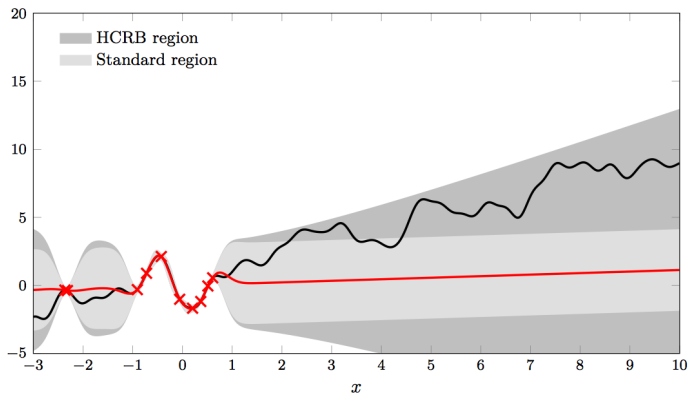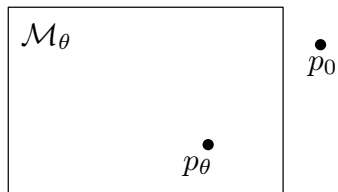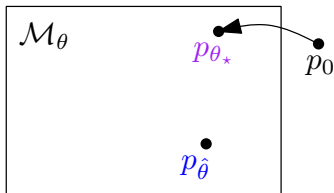
# Uncertainty region: underperformance



Figure : Standard region $C_\theta(x)$ with $\kappa = 3$ around prediction $\widehat{y}(x)$

# Misspecified case

# Seeking the best model



$$\boldsymbol{\theta}_\star = \arg\min_{\boldsymbol{\theta}} \ \Delta(\boldsymbol{\theta})$$

**Seeking the best model**

$$\mathcal{M}_\theta \qquad \overset{\bullet}{p_{\theta_\star}} \qquad \overset{\bullet}{p_0}$$

$$\overset{\bullet}{p_{\hat\theta}}$$

$$\boldsymbol{\theta}_\star = \underset{\boldsymbol{\theta}}{\arg\min}\ \Delta(\boldsymbol{\theta})$$

- How good is best model $\boldsymbol{\theta}_\star$?
- How far away is learned model $\widehat{\boldsymbol{\theta}}$?

# Uncertainty region: unknown properties



$$C(\mathbf{x})$$

$$y \qquad \hat{y}(\mathbf{x})$$

Credibility:

$$\Pr\left\{y \in C_{\widehat{\theta}}(\mathbf{x}) \mid \mathbf{x}, \mathcal{D}\right\} \quad ???$$

Confidence:

$$\Pr\left\{y \in C_{\widehat{\theta}}(\mathbf{x}) \mid \mathbf{x}\right\} \quad ???$$

dave.zachariah@it.uu.se

# Distribution-free approach
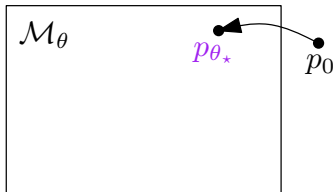
## Simple model class



Here: i.i.d. samples

$$\mathcal{M}_\theta = \Big\{ \, p_\theta(\mathbf{x}, \mathcal{D}) \, p_\theta(y | \mathbf{x}, \mathcal{D}) \, : \, \mu_\theta(\mathbf{x}) = \boldsymbol{\phi}^\top(\mathbf{x}) \boldsymbol{\theta} \, \Big\}$$

with $p$-dimensional feature vector $\boldsymbol{\phi}(\mathbf{x})$

# Defining the best model



$\mathcal{M}_\theta \qquad p_{\theta_\star} \qquad p_0$

Best model:

$$\boldsymbol{\theta}_\star = \underset{\boldsymbol{\theta}\,:\,\|\boldsymbol{\theta}\|_0 \leq k}{\arg\min}\ \Delta(\boldsymbol{\theta}) \quad \text{where} \quad \Delta(\boldsymbol{\theta}) = \widehat{\mathsf{E}}\Big[|y_i - \mu_\theta(\mathbf{x}_i)|^2\Big]$$

## Learning a model



Learned model:

$$\widehat{\boldsymbol{\theta}} = \arg\min_{\boldsymbol{\theta}} \ \sqrt{\Delta(\boldsymbol{\theta})} + \frac{1}{\sqrt{n}}\|\boldsymbol{\varphi} \odot \boldsymbol{\theta}\|_1,$$

with weight $\varphi_j = \|\widetilde{\boldsymbol{\phi}}_j\|_2/\sqrt{n}$ from feature $j$
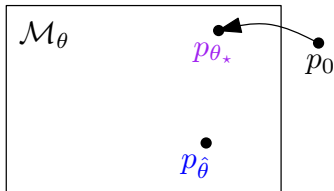
# Learning a model



### Guarantees

When $\varphi_j$ above a prediction error level of best model:

$$\widehat{\mathsf{E}}\Big[|\mu_{\theta_\star}(\mathbf{x}_i) - \mu_{\widehat{\theta}}(\mathbf{x}_i)|^2\Big] \leq \frac{2}{n}\|\boldsymbol{\varphi} \odot \boldsymbol{\theta}_\star\|_1^2 + 4\sqrt{\frac{\Delta(\boldsymbol{\theta}_\star)}{n}}\|\boldsymbol{\varphi} \odot \boldsymbol{\theta}_\star\|_1$$

# Learning a model



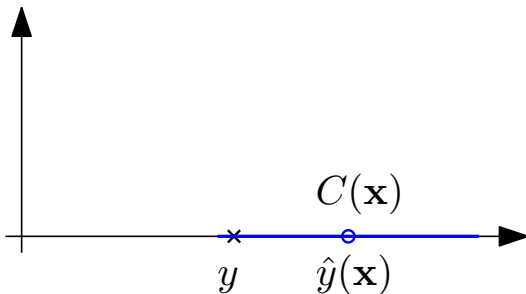$\mathcal{M}_\theta$    $p_{\theta_\star}$    $p_0$

$p_{\hat{\theta}}$

## Computational requirements

Recursive computation of $\widehat{\boldsymbol{\theta}}$ via matrix-inversion free updates with runtime $\mathcal{O}(p^2 n)$ and memory $\mathcal{O}(p^2)$
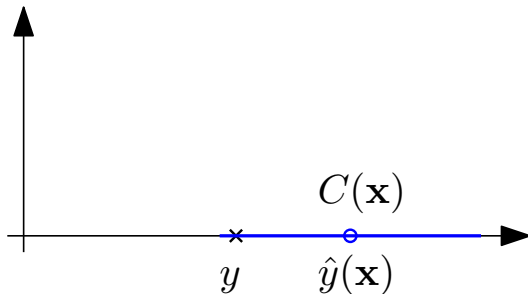
# Uncertainty region under misspecification



$$C(\mathbf{x})$$

$$y \qquad \hat{y}(\mathbf{x})$$

Constructing uncertainty region:

1. Randomly split $\mathcal{D}$ into $\mathcal{D}'$ and $\mathcal{D}''$
2. Learn $\widehat{\boldsymbol{\theta}}$ using $\mathcal{D}'$ and sort residuals $r_i = |y_i - \mu_{\widehat{\theta}}(\mathbf{x}_i)|$ from $\mathcal{D}''$
3. Let $\overline{r}_\kappa$ denote the $\lceil (n/2+1)\kappa \rceil$th smallest residual
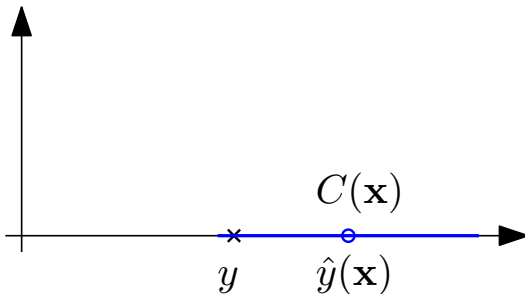
## Uncertainty region under misspecification



Construct uncertainty region

$$C_{\widehat{\theta}}(\mathbf{x}) = \left\{ y' \; : \; |y' - \mu_{\widehat{\theta}}| \leq \overline{r}_{\kappa} \right\}$$

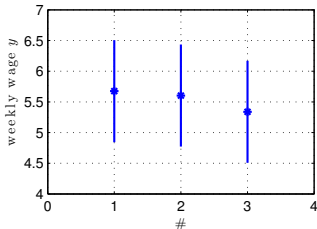# Uncertainty region under misspecification
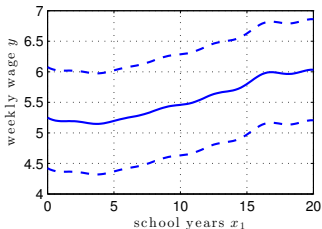


### Marginal confidence

When $(\mathbf{x}_i, y_i) \sim p_0(\mathbf{x}, y)$ independent,
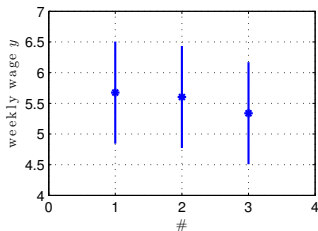
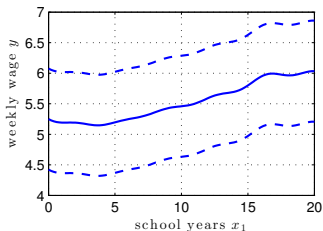$$\Pr\{\, y \in C_{\widehat{\theta}}(\mathbf{x}) \,\} \;\geq\; \kappa$$

# Empirical illustrations

## Example: US income data, 1930s cohort



- $\mathbf{x}$ years in school, marital status, ethnicity, region, etc.
- $y$ weekly wage [log-units]
- $n = 219\ 673$ samples
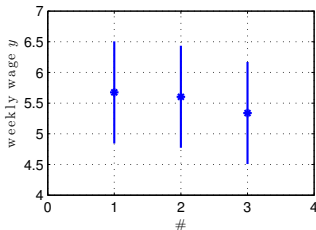- $\mathcal{M}_\theta$ where $\phi(\mathbf{x})$ is linear and wavelet-based
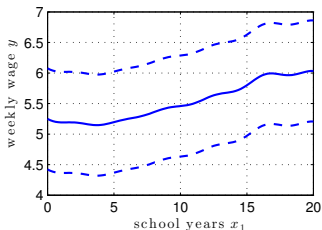
## Example: US income data, 1930s cohort



Left: Predicted wage vs. years in school.
Right: Predicted wage for individuals with 12 years of schooling, but differing backgrounds.
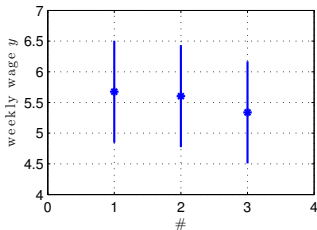
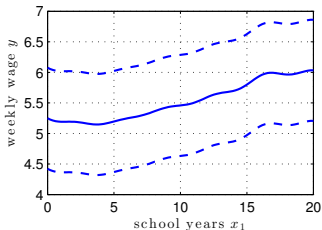# Example: US income data, 1930s cohort



Output dynamic range $y \in [-2.34, \ 10.53]$

$$\Rightarrow \widehat{\text{RMSE}} = 0.62$$

using $\bar{n} \sim 110 \times 10^3$ test individuals
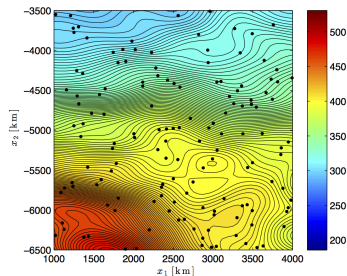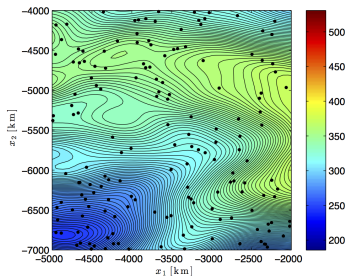
# Example: US income data, 1930s cohort



Output dynamic range $y \in [-2.34, \ 10.53]$

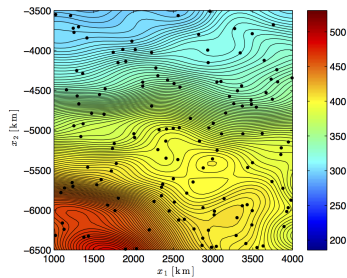$$\kappa = 0.90 \quad \Rightarrow \quad \widehat{\Pr\{y \in C_{\widehat{\theta}}\}} = 89.9\%$$

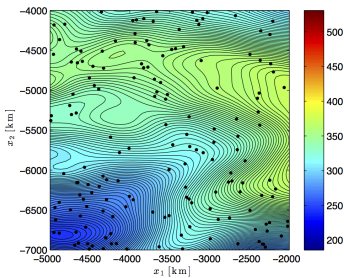using $\bar{n} \sim 110 \times 10^3$ test individuals

# Example: Global ozone data



- ▶ **x** spatial coordinates [km]
- ▶ $y$ ozone density [DU]
- ▶ $n = 17\ 340$ samples
- ▶ $\mathcal{M}_\theta$ where $\phi(\mathbf{x})$ is wavelet-based

# Example: Global ozone data



Output dynamic range $y \in [179, \ 542]$

$$\Rightarrow \widehat{\mathsf{RMSE}} = 6.74$$

using $\bar{n} \sim 156 \times 10^3$ test samples

# Example: Global ozone data



Output dynamic range $y \in [179, \, 542]$

$$\kappa = 0.90 \quad \Rightarrow \quad \widehat{\Pr\{y \in C_{\widehat{\theta}}\}} = 90.0\%$$

using $\bar{n} \sim 156 \times 10^3$ test samples

# Conclusions

# Conclusions

1. Subjectivist-Bayesian vs. Frequentist interpretations

# Conclusions

1. Subjectivist-Bayesian vs. Frequentist interpretations
   - Neither Bayes' rule nor prior are distinguishing features …
   - … rather belief vs. coverage

# Conclusions

1. Subjectivist-Bayesian vs. Frequentist interpretations
   - Neither Bayes' rule nor prior are distinguishing features ...
   - ... rather belief vs. coverage
2. Bayesian paradigm is powerful but possibly intoxicating

dave.zachariah@it.uu.se

# Conclusions

1. Subjectivist-Bayesian vs. Frequentist interpretations
   - Neither Bayes' rule nor prior are distinguishing features ...
   - ... rather belief vs. coverage
2. Bayesian paradigm is powerful but possibly intoxicating
   - Construct excellent decision rules ...
   - ... but also seriously misleading uncertainty regions

dave.zachariah@it.uu.se

# Conclusions

1. Subjectivist-Bayesian vs. Frequentist interpretations
   - Neither Bayes' rule nor prior are distinguishing features …
   - … rather belief vs. coverage
2. Bayesian paradigm is powerful but possibly intoxicating
   - Construct excellent decision rules …
   - … but also seriously misleading uncertainty regions
3. Distribution-free approach can tackle misspecification

# Conclusions

1. Subjectivist-Bayesian vs. Frequentist interpretations
   - Neither Bayes' rule nor prior are distinguishing features ...
   - ... rather belief vs. coverage
2. Bayesian paradigm is powerful but possibly intoxicating
   - Construct excellent decision rules ...
   - ... but also seriously misleading uncertainty regions
3. Distribution-free approach can tackle misspecification
   - Simple model class with meaningful best model
   - Performance guarantees and scalable in $n$ ...
   - ... uncertainty regions with valid coverage

dave.zachariah@it.uu.se