

Bayesian Approaches to Distribution Regression

Ho Chung Leon Law*, Dougal J. Sutherland*, Dino Sejdinovic, Seth Flaxman

ho.law@spc.ox.ac.uk, dougals@gatsby.ucl.ac.uk, dino.sejdinovic@stats.ox.ac.uk, s.flaxman@imperial.ac.uk

Distribution regression

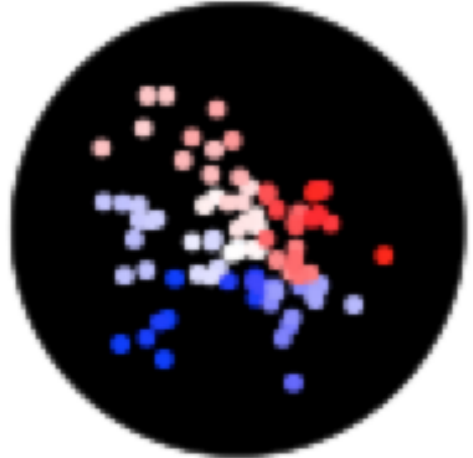
Inputs are sample sets from distributions [e.g. 7]:

$$y_i = f^*(\mathbb{P}_i) + \varepsilon \quad \mathbf{X}_i \stackrel{iid}{\sim} \mathbb{P}_i \quad \hat{y}_i = f(\mathbf{X}_i)$$

Examples:

\mathbb{P}_i	\mathbf{X}_i	y_i																				
Multivariate distribution [7]	Sample from distribution	Entropy of a projection																				
Demographics of county [2, 3]	<table border="1"> <thead> <tr> <th>AGEP</th> <th>SEX</th> <th>PINCP</th> <th>WKHP</th> </tr> </thead> <tbody> <tr> <td>45</td> <td>2</td> <td>14418.0</td> <td>30.0</td> </tr> <tr> <td>37</td> <td>1</td> <td>120152.0</td> <td>40.0</td> </tr> <tr> <td>63</td> <td>1</td> <td>11314.0</td> <td>NaN</td> </tr> <tr> <td>38</td> <td>2</td> <td>73092.0</td> <td>40.0</td> </tr> </tbody> </table>	AGEP	SEX	PINCP	WKHP	45	2	14418.0	30.0	37	1	120152.0	40.0	63	1	11314.0	NaN	38	2	73092.0	40.0	30.3% Clinton 9.9% Trump 59.8% none/other
AGEP	SEX	PINCP	WKHP																			
45	2	14418.0	30.0																			
37	1	120152.0	40.0																			
63	1	11314.0	NaN																			
38	2	73092.0	40.0																			

Galaxy velocities in a cluster [4, 5]



Cluster mass:
 $7 \times 10^{14} M_{\odot}$

Standard approach [e.g. 7]

• Kernel mean embeddings: kernel k , RKHS \mathcal{H} ,

$$\mu_{\mathbb{P}} = \mathbb{E}_{X \sim \mathbb{P}} [k(\cdot, X)] \in \mathcal{H}$$

• $\mu_{\mathbb{P}}$ fully characterizes \mathbb{P} for many k

• Empirical mean estimator:

$$\hat{\mu}_{\mathbb{P}} = \frac{1}{N} \sum_{j=1}^N k(\cdot, X^{(j)})$$

• Ridge regression from $\hat{\mu}_{\mathbb{P}}$ to y_i :

$$\begin{aligned} f(\hat{\mu}_{\text{test}}) &= \sum_{i=1}^n \alpha_i \langle \hat{\mu}_{\text{test}}, \hat{\mu}_i \rangle_{\mathcal{H}} \\ &= \sum_{i=1}^n \alpha_i \sum_{j=1}^{N_{\text{test}}} \sum_{j'=1}^{N_i} k(X_j^{\text{test}}, X_{j'}^i) \end{aligned}$$

• We make a landmark approximation (RBF net):

$$f(\hat{\mu}) = \beta^{\top} \hat{\mu}(\mathbf{u}) = \sum_{\ell=1}^s \beta_{\ell} \hat{\mu}(u_{\ell})$$

Sources of uncertainty

- Uncertainty about regression weights β
- Mean embeddings μ_i not seen exactly
- Should “trust” points with small N_i less

Bayesian linear regression

- Standard BLR model (similar to e.g. [3])
- Handles regression weight β uncertainty
- Assumes $\hat{\mu}_i$ known exactly
- Normal prior on regression weights $\beta \sim \mathcal{N}(0, \rho^2)$
- Observations $y_i | \mathbf{X}_i, \beta \sim \mathcal{N}(\beta^{\top} \hat{\mu}_i(\mathbf{u}), \sigma^2)$
- Gives normal $y_i | \mathbf{X}_i$ with hyperparameters σ, ρ

References

- [1] S. Flaxman, D. Sejdinovic, J. Cunningham and S. Filippi. ‘Bayesian Learning of Kernel Embeddings’. In: *UAI*. 2016. arXiv: 1603.02160.
- [2] S. Flaxman, D. J. Sutherland, Y.-X. Wang and Y.-W. Teh. *Understanding the 2016 US Presidential Election using ecological inference and distribution regression with census microdata*. 2016. arXiv: 1611.03787.
- [3] S. Flaxman, Y.-X. Wang and A. J. Smola. ‘Who Supported Obama in 2012?: Ecological inference through distribution regression’. In: *KDD*. ACM. 2015.
- [4] M. Ntampaka, H. Trac, D. J. Sutherland, N. Battaglia, B. Póczos and J. Schneider. ‘A Machine Learning Approach for Dynamical Mass Measurements of Galaxy Clusters’. In: *The Astrophysical Journal* 803.2 (2015). arXiv: 1410.0686.
- [5] M. Ntampaka, H. Trac, D. J. Sutherland, S. Fromenteau, B. Póczos and J. Schneider. ‘Dynamical Mass Measurements of Contaminated Galaxy Clusters Using Machine Learning’. In: *The Astrophysical Journal* 831.2 (2016), p. 135.
- [6] R. Rothe, R. Timofte and L. V. Gool. ‘Deep expectation of real and apparent age from a single image without facial landmarks’. In: *IJCV* (July 2016).
- [7] Z. Szábo, B. K. Sriperumbudur, B. Póczos and A. Gretton. ‘Learning Theory for Distribution Regression’. In: *JMLR* 17.152 (2016), pp. 1–40. arXiv: 1411.2066.

Shrinkage

- New model
- Handles $\hat{\mu}_i$ uncertainty
- Point estimate for β to retain conjugacy
- Uses Bayesian nonparametric model for $\hat{\mu}_i$ [1]:
- Prior: $\mu_i \sim \mathcal{GP}(m_0, \eta r(\cdot, \cdot))$
- Likelihood: “observed” at points \mathbf{u} , CLT:

$$\hat{\mu}_i(\mathbf{u}) | \mu_i(\mathbf{u}) \sim \mathcal{N}(\mu_i(\mathbf{u}), \Sigma_i)$$
- Closed-form GP posterior for $\hat{\mu}_i | \mathbf{X}_i$
- Similar to James-Stein shrinkage

• Observations $y_i | \mu_i, f \sim \mathcal{N}(\langle f, \mu_i \rangle_{\mathcal{H}}, \sigma^2)$

• Predictive: $y_i | \mathbf{X}_i, \beta, \mathbf{z} \sim \mathcal{N}(\xi_i^{\beta}, \nu_i^{\beta})$

$$\xi_i^{\beta} = \beta^{\top} R_{\mathbf{z}} \left(R + \frac{\Sigma_i}{N_i} \right)^{-1} (\hat{\mu}_i - m_0) + \beta^{\top} m_0$$

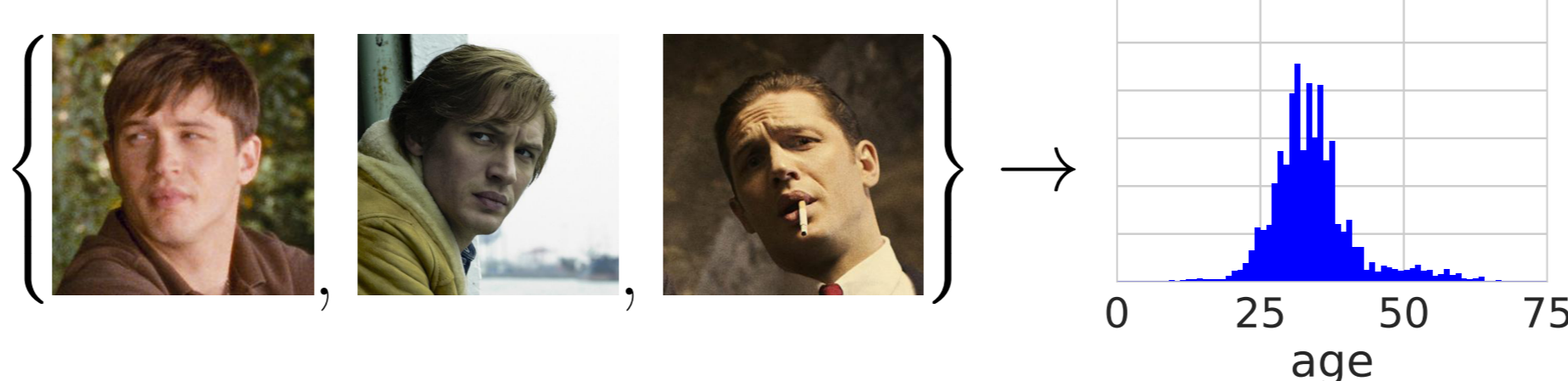
$$\nu_i^{\beta} = \beta^{\top} \left(R_{\mathbf{z}\mathbf{z}} - R_{\mathbf{z}} \left(R + \frac{\Sigma_i}{N_i} \right)^{-1} R_{\mathbf{z}}^{\top} \right) \beta + \sigma^2$$

- Point estimates for β, σ , maybe kernel params. . .
- Optimise with backprop (TensorFlow)

Bayesian Distribution Regression

- Full uncertainty model (BDR)
- Shrinkage posterior for $\hat{\mu}_i$
- BLR-like posterior for weights α
- Not fully conjugate
- MCMC for inference about α (Stan)

Age prediction from images



- IMDb database from [6]
- *Very* noisy labels in the dataset
- Group pictures of actors, predict mean age
- Features: last hidden layer from [6]’s CNN
- Lots of variation in N_i :

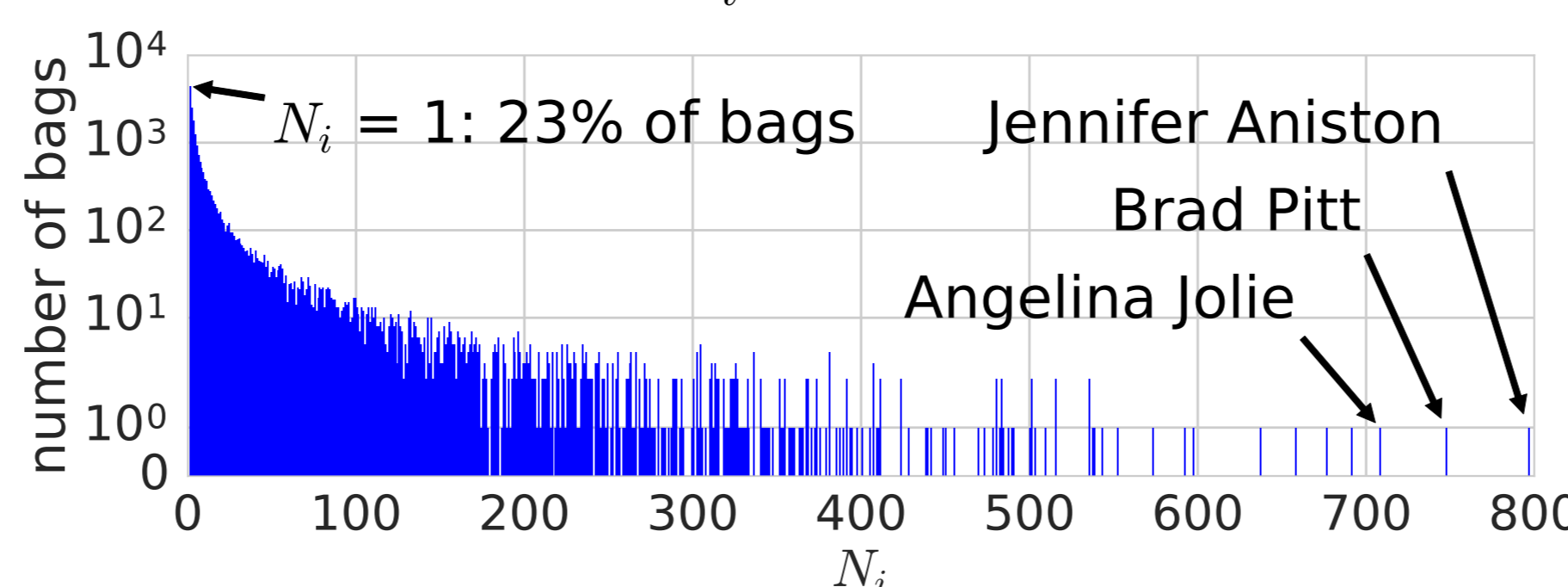


Figure: Histogram of N_i .

• Shrinkage really helps!

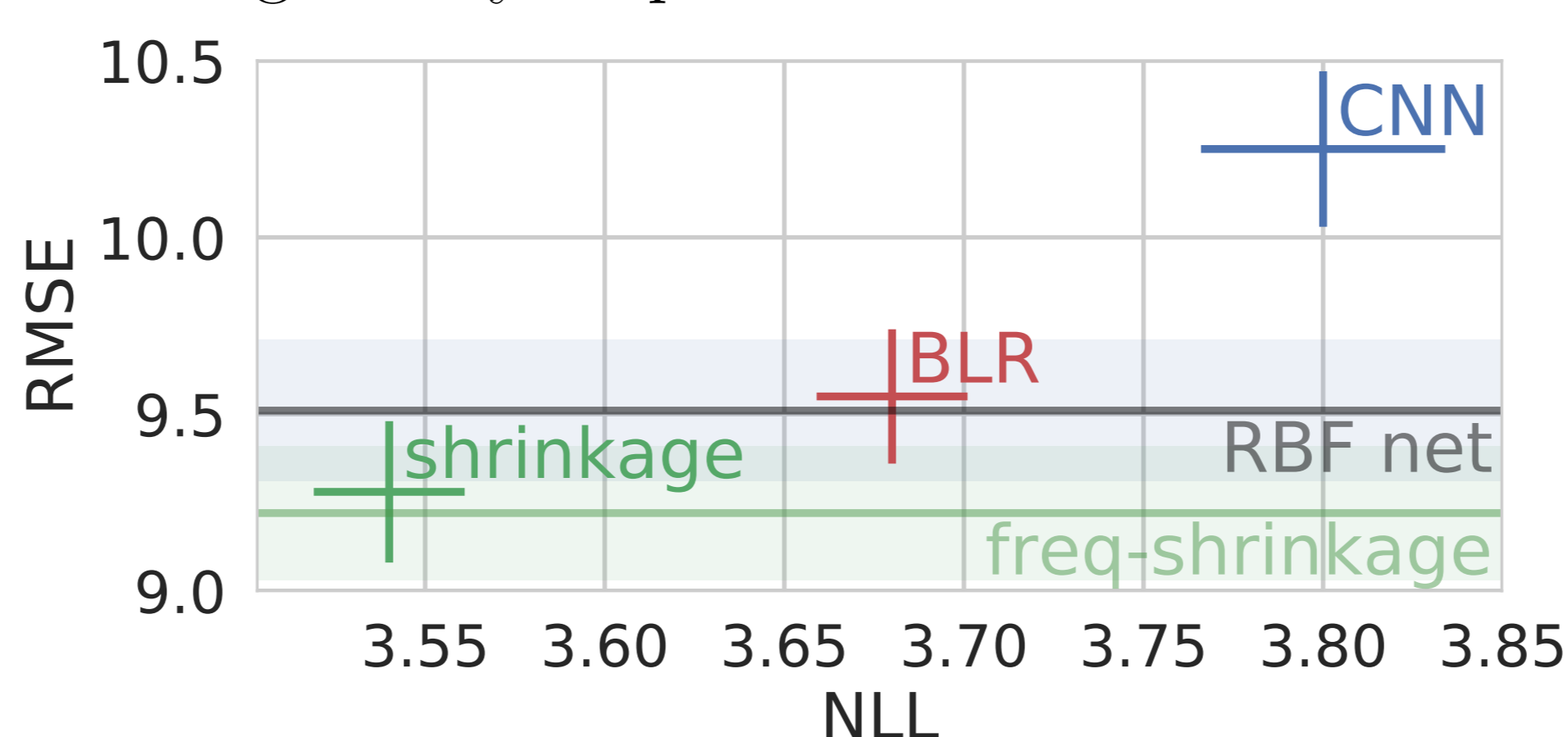


Figure: Results across ten data splits (means and standard deviations). RBF net, freq-shrinkage are tuned for RMSE, other methods for NLL. CNN takes the mean of the predictive distributions of [6] for each point in the bag.

Synthetic data

- Toy problem to examine input uncertainty:
- Labels y_i uniform over $[4, 8]$
- 5d data points: $[x_j^i]_{\ell} | y_i \stackrel{iid}{\sim} \frac{1}{y_i} \Gamma\left(\frac{y_i}{2}, \frac{1}{2}\right)$
- $(s_5, 25, 25, 100 - s_5)\%$ of bags respectively have $N_i = (5, 20, 100, 1000)$

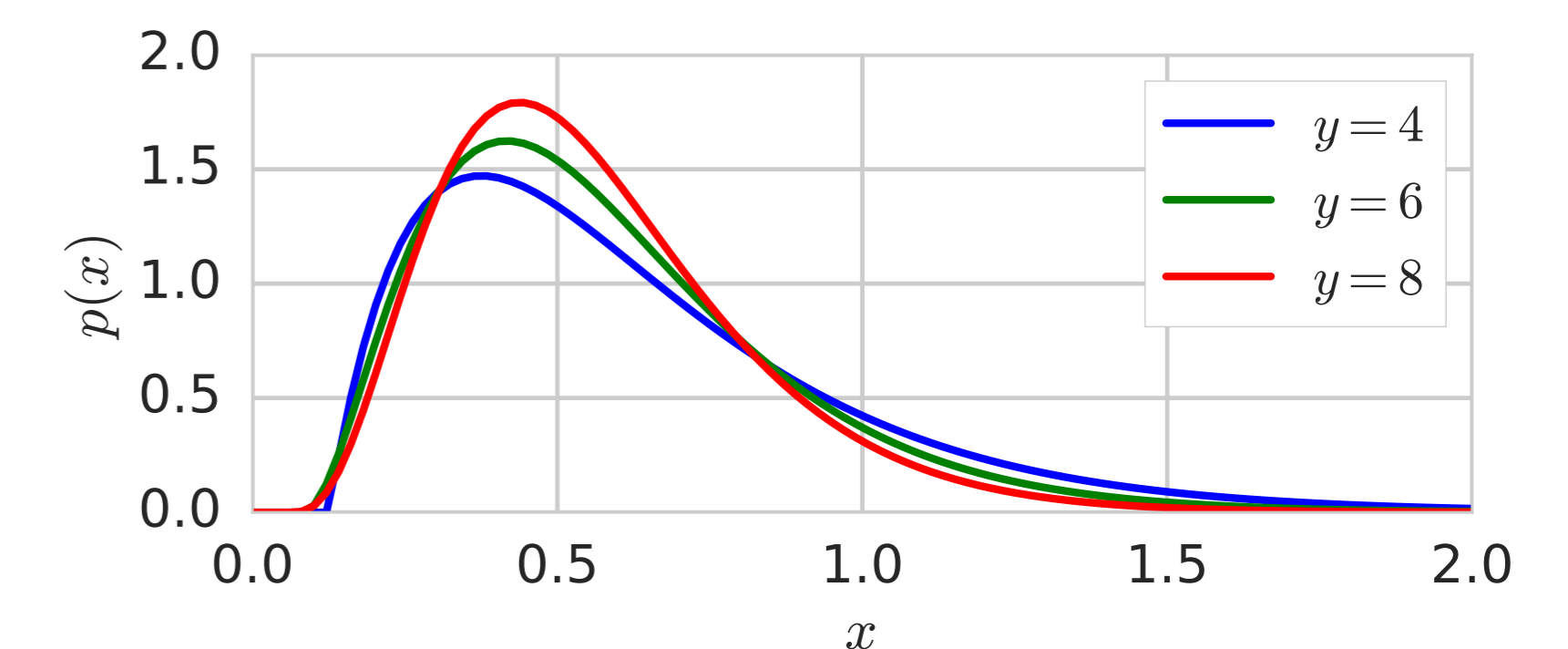


Figure: Density of x_j^i for different y_i .

• BDR \approx shrinkage < BLR in NLL, MSE:

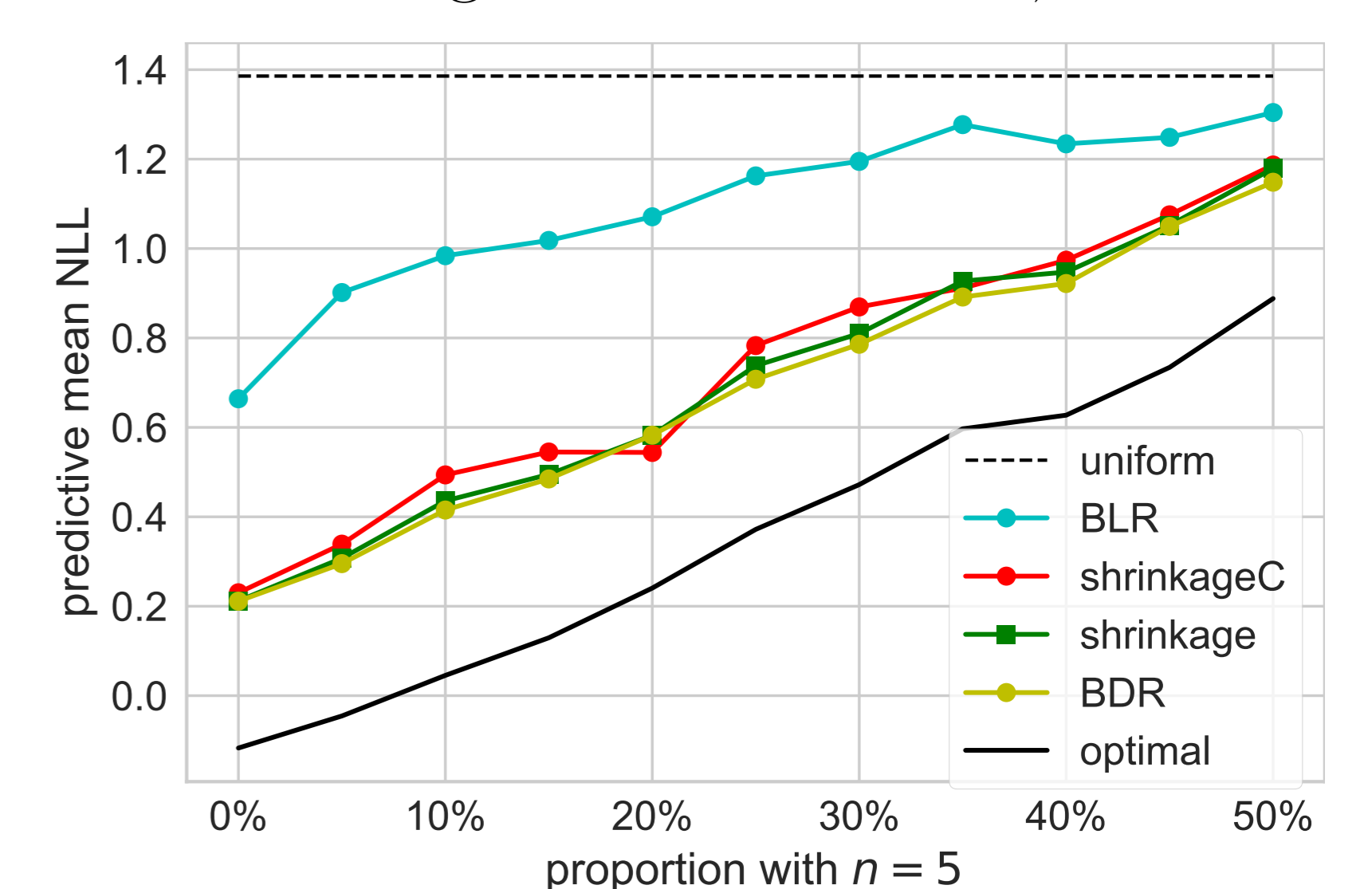


Figure: Negative log predictive likelihoods.

• Illustrative behaviour of different kinds of bags:

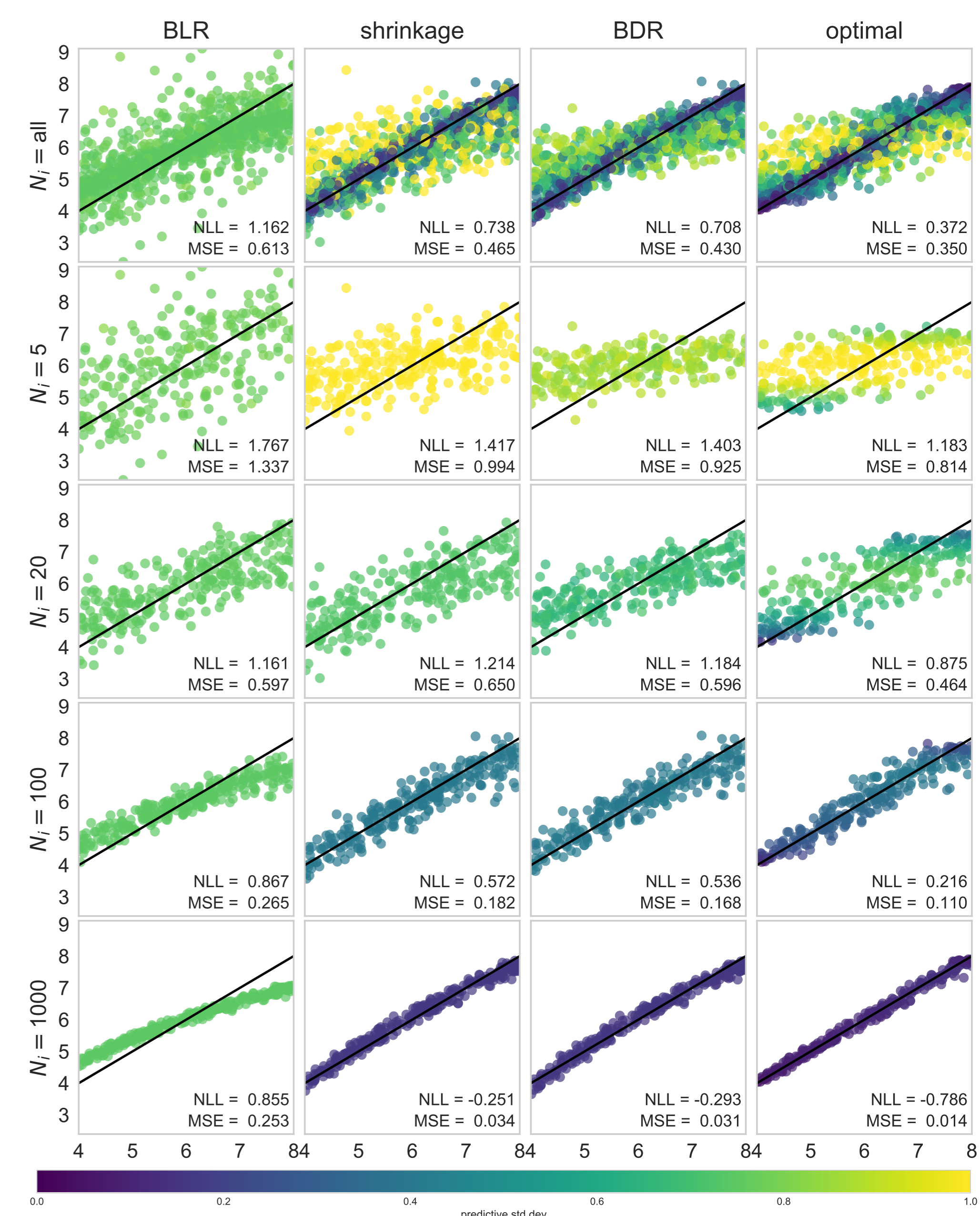


Figure: Bags are dots: horizontal position is true y_i , vertical is predictive mean, colour is predictive std. $s_5 = 25$.

- Similar experiment with $N_i = 1000$, label noise:
- BDR \approx BLR < shrinkage in NLL, MSE
- BDR able to take advantage of both settings

Takeaway

- Two sources of uncertainty in distribution regression: inputs and model
- BLR handles model uncertainty
- Shrinkage method handles input uncertainty based on bag size
- Full BDR handles both, but needs MCMC